



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Enhancing toxicity prediction of MLinvitroTox: Prioritizing unidentified compounds in environmental samples based on hazard assessment

Master Thesis

Robin Bosshard, 16-915-399

October 16, 2023

Supervisors: Prof. Dr. Fernando Perez Cruz, Dr. Eliza Harris, Lili Gasser (SDSC)
Dr. Kasia Arturi (Eawag)

Department of Computer Science, ETH Zürich

Abstract

This thesis enhances the capabilities of the MLinvitroTox framework, which is designed to forecast the toxicity of unidentified chemical compounds utilizing High-Resolution Mass Spectrometry (HRMS/MS) data. The framework aims to maximize the detection rate of most hazardous compounds within environmental samples, all while minimizing the occurrence of false alarms and channeling resources and efforts into the labor-intensive task of identifying and quantifying compounds that possess the highest potential for causing harm. This approach stands in contrast to standard nontarget screening HRMS workflows, which typically prioritize compounds that are detected most frequently and with highest intensities in the mass-to-charge ratio of ions in a sample. We employed hazard-driven machine learning models based on molecular fingerprints derived from chemical structure and utilized *in vitro* toxicity data from ToxCast/Tox21. We have developed pytcp1, a Python processing pipeline that is applicable to the latest toxicity data. We have leveraged datasets spanning various assay endpoints that encompass a wide range of toxicity aspects. The XGBoost classifier achieves a median F1-score of 0.65 across all target assay endpoints when predicting binary toxicity based on molecular fingerprints from known chemical structures. The models also demonstrate effectiveness in predictivity when validated on SIRIUS predicted molecular fingerprints from MassBank spectra. Furthermore, a web app was created to facilitate interaction with the data and the MLinvitroTox framework.

Acknowledgments

First and foremost, I would like to thank Prof. Dr. Fernando Perez Cruz from the Swiss Data Science Center (SDSC) for granting me the opportunity to work on this fascinating project. His support has been invaluable.

I would like to express my sincere gratitude to my supervisor Dr. Kasia Arturi from Swiss Federal Institute of Aquatic Science and Technology (Eawag) and my supervisors Dr. Eliza Harris, Lili Gasser from SDSC for their numerous discussions, patience and valuable insights. Without their help, this thesis would not have been achievable.

Additionally, I would like to acknowledge Prof. Dr. Juliane Hollender from Eawag for her support throughout the project and for the enlightening experience of visiting the Eawag labs.

Furthermore, my gratitude goes out to Jason Brown, Feshuk Madison, and Katie Paul Friedman from U.S. EPA for their participation in discussions concerning the technical aspects of the tcpl pipeline and the ToxCast database.

Lastly, I extend a special thank you to my family and friends for their unconditional support throughout my academic journey.

Contents

Contents	iii
1 Introduction	1
1.1 The Challenge of Environmental Pollution	1
1.2 The Imperative for Prioritization and Toxicity Assessment	3
1.3 Unlocking the Potential of High-Throughput Screening and Machine Learning in Toxicity Prediction	4
1.4 MLinvitroTox: A Novel Approach	5
1.5 Objectives and Significance	6
1.6 Thesis Structure	6
2 Background	7
2.1 Toxicity Testing: From In Vitro Assays and Molecular Fingerprints to Predictive Models and Beyond	7
2.2 Chemical Target Toxicity vs. Cytotoxicity	10
3 Related work	12
4 Material and Methods	14
4.1 Toxicity Data and Processing	14
4.1.1 ToxCast invitroDB v4.1	14
4.1.2 tcpl v3.0	14
4.1.3 Concentration-Response Series	15
4.1.4 Efficacy Cutoff	17
4.1.5 tcplFit2	17
4.1.6 Curve Fitting	17
4.1.7 Hit Calling	20
4.1.8 Flagging	21
4.2 New Toxicity Pipeline Implementation: pytcpl	21
4.2.1 Introduction	21

4.2.2	Setup step	21
4.2.3	Main step	23
4.2.4	Post-Processing step	24
4.2.5	Curve Surfer	25
4.3	Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline	26
4.3.1	Training	27
4.3.2	Evaluation	29
4.3.3	Application	31
5	Results	33
5.1	Binary Classification	33
5.1.1	Evaluation Metrics	33
5.1.2	Performance	37
5.1.3	Feature Importance	45
6	Discussion	46
7	Conclusion	48
Bibliography		49
A	Appendix	53
A.1	Variability in the Tested Concentration Across Assay Endpoints and Compounds	53
A.2	ToxCast Assay Sources	54

Chapter 1

Introduction

1.1 The Challenge of Environmental Pollution

Over the past few decades, the upsurge in environmental pollution by chemical compounds has been driven by industrial processes, agricultural methods, consumerism and various other contributing factors. Although these chemicals are integral for many products and have the potential to improve the comfort of modern society, they can also pose risks and adversely affect both human health and the environment, either acutely or chronically. Toxic substances threaten wildlife but also make air, soil, drinking water and food supply less safe.

Nations worldwide maintain comprehensive chemical regulations¹, however, it is anticipated that global chemicals production will double by 2030 [1]. Moreover, the widespread utilization of chemicals, including their use in consumer goods, is expected to expand further. Even though there are over 275 million known chemical compounds registered by the Chemical Abstracts Service [2], merely a tiny fraction of them undergo close monitoring via target analytical approaches and even less is known about their toxicity profiles and negative health effects on organisms. Refer to Table 1.1 for an overview of omnipresent water pollutants.

In light of the rapidly evolving chemical landscape, there is an increasing demand for future-proof, robust measurement and modeling methods. These methods are crucial for evaluating the toxicity of chemicals, facilitating informed risk-based decision-making even when data on hazards and exposures are limited. It is worth noting that the need for adaptable approaches in chemical safety and sustainability efforts must also prioritize cost-efficiency and gain widespread acceptance among regulatory bodies, industry stakeholders, and the general public. For instance, the EU has introduced the 8th Environment Action Programme, as outlined in its European Green Deal [4], to provide direction for European environmental policy until the year 2030. This

¹For instance, REACH, short for Registration, Evaluation, Authorisation, and Restriction of Chemicals, is an EU regulation aimed at improving chemical safety and allocating risk management responsibilities to companies operating in various sectors.

1.1. The Challenge of Environmental Pollution

Table 1.1: Examples of ubiquitous water pollutants. Table 2 adapted from [3].

Origin/Usage	Class	Examples	Related Issues
Industrial Chemicals	Solvents	Tetrachloromethane	Hepatotoxicity
	Intermediates	Methyl-t-butylether	Drinking-water-quality
	Petrochemicals	BTEX	Cancer
Industrial Products	Additives	Phthalates	Endocrine disruptors
	Lubricants	PCBs	Biomagnification
	Flame Retardants	PBDEs	
Consumer Products	Detergents	Nonylphenol ethoxylates	Endocrine effects
	Pharmaceuticals	Antibiotics	Bacterial resistance
	Hormones	Ethynodiol diacetate	Feminization of fish
Biocides	Pesticides	DDT	Toxic effects and persistent metabolites
Natural Chemicals	Heavy Metals	Lead, mercury	Organ damage
	Inorganics	Arsenic, fluoride	Drinking-water-quality
	Taste and Odor	Geosmin	
	Hormones	Estradiol	Feminization of fish
Disinfection & Oxidation	Disinfection by-products	Haloacetic acids, Bromate	Drinking-water-quality
Transformation Products	Metabolites from all above	Metabolites of perfluorinated compounds	Bioaccumulation

program reinforces the EU's ambitious goal of sustainable living within planetary limits, with a forward-looking vision that extends to 2050. Central to this vision is a zero-pollution commitment, encompassing air, water, and soil quality. In 2021, the European Commission introduced a sustainability-focused chemicals strategy [5], which aligns with the EU's zero-pollution ambition. This strategy not only enables the evaluation of the safety and sustainability of emerging compounds but also aims to reduce existing concerning substances, such as *per- and polyfluoroalkyl substances* (PFAS), through substitution or phasing out wherever feasible. In parallel, the U.S. *Environmental Protection Agency* (EPA) shares a similar scientific consensus and is at the forefront of assessing the potential impacts of chemicals on human health and the environment. Leveraging advanced toxicological methods, EPA actively promotes risk reduction efforts through its own Chemical Safety for Sustainability National Research Program. This program builds upon the achievements of research initiatives like *ToxCast/Tox21*² and the Endocrine Disruptor Screening Program in the 21st Century (EDSP21), demonstrating a commitment to advancing chemical safety on a global scale.

²<https://www.epa.gov/chemical-research/exploring-toxcast-data>

1.2 The Imperative for Prioritization and Toxicity Assessment

Modern analytical techniques, including *high resolution mass spectrometry (HRMS/MS)*, are gaining significance across various domains such as metabolomics, drug discovery, environmental science and toxicology [6].

In environmental monitoring, the application of nontarget HRMS/MS has notably improved the capacity to detect possibly thousands of contaminants in a single sample. The instrument generates complex spectra that provide information about the masses and fragmentation patterns of compounds present within the sample as illustrated in Figure 1.1. Often, only a minority of these molecules can be definitively identified, while the majority remains unidentified, resulting in their classification into two categories:

- **Target compounds** are substances that researchers intentionally seek to identify and quantify in a given sample due to their known importance or relevance. These compounds have well-documented chemical structures and properties, and their identities have been confirmed using various analytical techniques. Target compounds are linked to existing databases or reference spectra whenever available, making it easier to access information about their toxicity and other relevant characteristics. The selection of target compounds is a common practice in analytical chemistry and environmental monitoring and is crucial for focusing research efforts on specific compounds of interest.
- **Non-target compounds**, on the other hand, are substances that are not intentionally selected as the primary focus of analysis. They are substances detected in the sample but lack definitive characterization in terms of their chemical identity, structure, or properties, including their potential toxicity. These compounds are observed as peaks or features in spectral data, but their specific chemical attributes remain unknown. Identifying non-target compounds requires further investigation, which can be a resource-intensive process. Therefore, prioritizing the examination of non-target compounds becomes essential to efficiently allocate resources and gain a comprehensive understanding of a sample's chemical composition.

When it comes to prioritizing unidentified compounds for further in-depth testing and identification, the standard approach has been to rely on signal intensity in the chromatographic peaks. However, this approach tends to fall short in delivering an accurate assessment of environmental exposures because the signal intensity may not relate proportionally to the compound's concentration in the sample. Furthermore this approach overlooks the toxicological factors essential for prioritizing compounds with concerns related to environmental hazards and pollution. As a result, substances with the potential for severe ecological consequences, such as endocrine-disrupting compounds, often go undetected because of their low frequency, even though they exhibit high levels of toxicity. Hence, a pressing need exists for alternative approaches to prioritize unidentified nontarget HRMS/MS signals based on their hazard potential. By incorporating relevant toxicity factors into the equation:

1.3. Unlocking the Potential of High-Throughput Screening and Machine Learning in Toxicity Prediction

$$\text{Risk} = \text{Hazard} \times \text{Exposure} \quad (1.1)$$

we augment the capacity to make well-informed decisions when evaluating the environmental risk associated with chemicals.

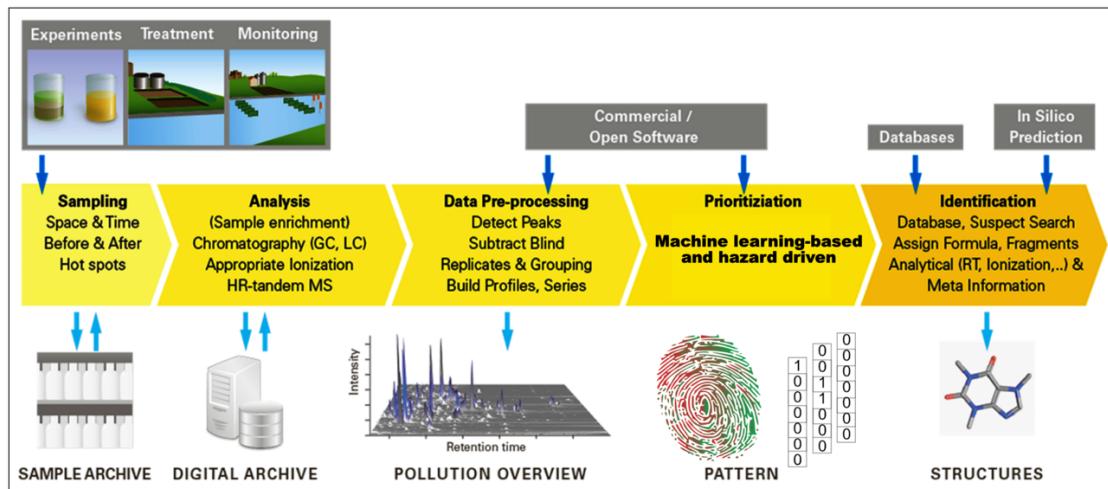


Figure 1.1: Schematic of the workflow used for nontarget HRMS/MS screening of environmental samples, featuring a customized prioritization step. Adapted from Figure 1 in the original source [7].

1.3 Unlocking the Potential of High-Throughput Screening and Machine Learning in Toxicity Prediction

In the past few years, the use of machine learning methods has emerged as a transformative force in the field of *in vitro* toxicology, particularly in the realm of high-throughput toxicity prediction. *High-throughput screening (HTS)* has revolutionized the way toxicity is assessed by allowing thousands of *in vitro* bioassays to be conducted efficiently. This high-throughput approach, coupled with advancements in robotics and automated analysis, has generated large volumes of toxicity data, paving the way for more comprehensive assessments of chemical compounds. Alongside the rise of machine learning, this advancement has facilitated the creation of predictive models, known as Quantitative Structure-Activity Relationship (QSAR) models. In the context of this research, we are specifically concerned with *Quantitative Structure-Toxicity Relationship (QSTR)* models. These models are capable of forecasting compound toxicity based on their physico-chemical properties or molecular descriptors [8]. As they are trained on extensive datasets containing toxicity information, these models can learn the underlying patterns and relationships between chemical structures and target toxicity. With this capability, they can predict the toxicity of new compounds, even when these substances themselves have not undergone laboratory testing. This approach holds the potential to

substantially decrease the time and expenses linked to initial toxicity pre-assessment, and it plays a pivotal role in determining which compounds should undergo more in-depth testing.



(a) A robot arm retrieves assay plates from incubators and places them at compound transfer stations or hands them off to another arm that services liquid dispensers or plate readers. Efforts in the automation, miniaturization and the readout technologies have enabled the growth of HTS. Image obtained from [9].

(b) Modern microtitre assay plates consist of multiples of 96 wells, which are either prepared in the lab or acquired commercially from stock plates. These wells are filled with a dilution solvent, such as *Dimethylsulfoxide (DMSO)*, along with the chemical compounds intended for analysis. Image obtained from [10].

Figure 1.2: High-Throughput Screening (HTS)

1.4 MLinvitroTox: A Novel Approach

In response to the pressing need for a more hazard-driven and inclusive assessment of environmental contaminants, Arturi *et al.* introduced *MLinvitroTox* [11], an innovative machine learning framework. This framework is part of a broader pipeline named *EXPECTmine*, which incorporates the complementary exposure aspect within the risk assessment process. The primary objective of this thesis is to collaborate with the authors to further enhance and advance this framework. *MLinvitroTox* leverages molecular fingerprints extracted from fragmentation spectra, marking a significant change in how the toxicity of the myriad unidentified HRMS/MS features is forecasted. *MLinvitroTox* follows a similar training approach as traditional QSTR models, using supervised classification models trained with molecular fingerprints derived from chemical structures. However, during the application phase, the input to the machine learning model consists of molecular fingerprints generated from experimentally measured mass-spectrometry fragmentation spectra using *SIRIUS+CSI:FingerID* [12]. *SIRIUS* is a software package for annotating small molecules from nontarget HRMS/MS data, while *CSI:FingerID* is a machine-learning tool employed by *SIRIUS* to predict molecular fingerprints from fragmentation spectra. Utilizing streamlined machine learning methodologies, *MLinvitroTox* forecasts chemical toxicity for a wide range of compounds. This analysis covers more than 300 target-specific assay endpoints, drawing data from ToxCast/Tox21

--- 1.5. Objectives and Significance

datasets. Subsequently, the toxicity predictions generated by the framework are employed to prioritize compounds, with the flexibility to emphasize specific aspects of toxicity profiles.

1.5 Objectives and Significance

The main objective of this thesis is to contribute to the development of an efficient MLinvitroTox framework for predicting compound toxicity across multiple endpoints. The goal is to enhance the integration of MLinvitroTox by creating an automated pipeline in the Python programming language. This pipeline is designed to efficiently address the inherent complexities associated with modeling and processing heterogeneous datasets. In this context, the primary focus is on elevating the quality of curating and preparing toxicological data, with a particular emphasis on streamlining the entire process. This process begins with raw concentration-response series data and ultimately leads to the generation of conclusive toxicity predictions. The ultimate output is expected to comprise *toxicity fingerprints* that encapsulate the predicted toxicity from HRMS/MS environmental samples for the relevant endpoints of interest. These generated toxicity fingerprints will offer crucial insights for the prioritization process, aiding in the identification of the most hazardous compounds present in environmental samples.

One notable constraint of the existing framework lies in its binary *hitcall* when predicting the toxicity of specific endpoints. It categorizes compounds as either toxic or non-toxic without accounting for variations in toxicity severity. In the long term, it is essential to adopt a more refined approach that can capture the nuanced continuum of toxicity. This thesis endeavors to overcome this limitation by developing a pipeline capable of forecasting toxicity across numerous endpoints, employing continuous hitcalls.

1.6 Thesis Structure

In the course of progressing through the subsequent chapters, insights will be provided into the materials and methods employed, focusing on the technical intricacies involved in the preparation of ToxCast/Tox21 toxicity data and their transformation into suitable inputs for the machine learning pipeline. This foundational work will establish the basis for the upcoming chapters, which will showcase the potential of MLinvitroTox. Furthermore, the framework's effectiveness is demonstrated through the validation of real-world mass spectral data from *MassBank* [13], and the examination of the implications of this research is carried out.

Chapter 2

Background

This chapter is vital for understanding the following sections of this thesis as it provides some foundational background information in toxicity testing.

2.1 Toxicity Testing: From In Vitro Assays and Molecular Fingerprints to Predictive Models and Beyond

With the ever-growing amount of chemical compounds entering the environment, traditional experimentation methods face limitations concerning cost and time constraints. Additionally, ethical concerns arise regarding the use of animal trials in *in vivo* experiments.

In 2007, the *U.S. National Academy of Sciences* introduced a visionary perspective and published a landmark report, titled as *Toxicity Testing in the 21st Century: Vision and Strategy*. This report promoted a transition from conventional, resource-consuming animal-based *in vivo* tests to efficient high-throughput *in vitro* pathway assays on cells. This transition paved the way for the realm of HTS, where a multitude of *in vitro* bioassays can be executed, complementing and improving chemical screening. This transformation is made possible by advancements in robotics, data processing, and automated analysis. As a result, this synergy has led to the generation of extensive toxicity datasets like ToxCast/Tox21.

HTS datasets, including ToxCast and other sources, have opened the door to promising applications of machine learning in predictive computational toxicology. These predictive models can be developed to screen environmental samples with limited availability of toxicity data, allowing for the prioritization of further testing efforts. Such models often forecast toxicity using QSTRs, which are based on descriptors encoding chemical structures like molecular fingerprints. 1D-Molecular fingerprints encode compound molecules as fixed-length binary vectors, denoting the presence (1) or absence (0) of specific substructures or functional groups, visualized in 2.1. Typically, fingerprints use *SMARTS* strings, as an extension of *SMILES* strings, to encode the underlying

2.1. Toxicity Testing: From In Vitro Assays and Molecular Fingerprints to Predictive Models and Beyond



Figure 2.1: Schematic of a molecular fingerprint for a fictional chemical. Each bit position accounts for the presence or absence of a specific structural fragment. Bit positions are set on (set to 1, gray) if the substructure is present in a molecule, or set off (set to 0, white) if it is absent. Figure 1 adapted from [14].

substructural patterns within molecules. While SMILES is a widely accepted notation system for representing chemical structures, there can be variations in how different sources generate SMILES strings. These variations can include the presence or absence of hydrogens, different ways of representing aromatic rings, variations in chemotypes, and tautomer representations. SMILES strings for the same chemical can differ due to these variations. To ensure consistency and deterministic computation, chemists often normalize SMILES before use, ensuring adherence to a common set of rules for computational analysis.

Once generated, molecular fingerprints can be used for various cheminformatics tasks. For example, they can be compared to identify structurally similar compounds, used as input for machine learning models to predict properties or activities.

To go one step further, SIRIUS employs CSI:FingerID, a method that directly predicts various fingerprint types from HRMS/MS fragmentation spectra. CSI:FingerID utilizes machine learning techniques, encompassing linear Support Vector Machines and Deep Learning, to predict an array of fingerprints, including CDK Substructure, PubChem CACTVS, Klekota-Roth, FP3, MACCS, ECFP2, and ECFP4 fingerprints. Different fingerprint types used in cheminformatics capture varying aspects of chemical compounds, such as molecular features, granularity, size, encoding method and domain-specific requirements.

2.1. Toxicity Testing: From In Vitro Assays and Molecular Fingerprints to Predictive Models and Beyond

The utilization of molecular fingerprints for *in vitro* toxicity prediction is based on the assumption that molecular toxic effects result from interactions between distinct chemical components and receptors during a *molecular initiating event (MIE)*. On a larger biological scale, the MIE can set a sequential chain of causally linked *key events (KE)* in motion. This occurs at different levels of biological organization from within cells to potentially culminating in an *adverse outcome pathway (AOP)* at the organ or organism level, as depicted in Figure 2.2. The mechanistic information captured in AOPs reveal how chemicals or other stressors cause harm, offering insights into disrupted biological processes, potential intervention points but also guide regulatory decisions on next generation risk assessment and toxicity testing. The AOP framework is an analytical construct that allows an activity mapping from the presence or absence of certain molecular substructures encoded in chemical descriptors to the target mechanistic toxicity. Finally, when monitoring disruptions in toxicity pathways, physiologically based pharmacokinetic (PBPK) models can be leveraged to extrapolate *in vitro* findings to human blood and tissue concentrations [15].

It is important to emphasize that the predictions from HTS bioassays portray molecular toxicity events only at a cellular level, and their translation to adverse outcomes at higher organism levels is not necessarily guaranteed. As the scale shifts from the cellular to the organism level, the confidence in these relationships may decrease.

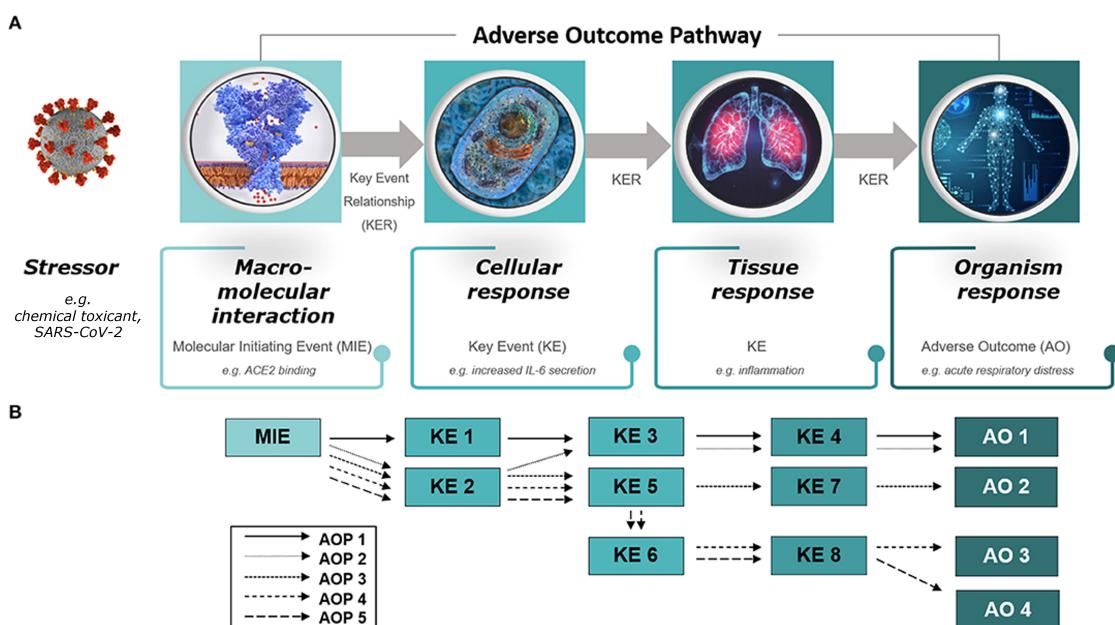


Figure 2.2: Diagram of (A) an adverse outcome pathway (AOP) and (B) an AOP network. (A) An AOP starts with a molecular initiating event (MIE), followed by a series of key events (KEs) on different levels of biological organization (cellular, tissue, organ) and ends with an adverse outcome (AO) in an organism. The stressor is not part of the AOP itself. Figure adapted from [16]

2.2 Chemical Target Toxicity vs. Cytotoxicity

Consider a hypothetical scenario in which a chemical undergoes testing in a bioassay that assesses toxicity by measuring the activation of a reporter gene within a cell. The reporter gene encodes a detectable protein, and its activation is triggered by the chemical binding to a specific receptor, the key focus of the assay endpoint. While it might seem logical that an increase in chemical concentration would result in a higher chemical toxicity signal, this assumption does not hold true in general. At elevated concentrations, the chemical can become *cytotoxic*, causing harm to the cells and ultimately leading to cell death. Consequently, this can lead to a decrease in the activation of the reporter gene and a subsequent reduction in the signal, indicating a decrease in bioactivity. Alternatively, in a last ditch effort to survive, a phenomenon known as the *cytotoxicity burst* [17] may occur. During this phenomenon, all cellular mechanisms, including the target mechanism, may be activated, leading to a burst in bioactivity. For a visual representation, consult Figure 2.3. Considering this situation, chemical toxicity can manifest in various forms, categorizing into two primary groups [17]:

- **Specific toxicity** occurs when a chemical interacts with and interferes with a specific biomolecular target or pathway, manifesting as effects like receptor agonism/antagonism or enzyme activation/inhibition. This thesis primarily focuses on specific toxicity, which is often the desired signal to detect in a target assay endpoint. However, it is essential to recognize that data processing must also take into account the following:
- **Non-specific toxicity (Cytotoxicity and cell stress)** involve broad disruptions of the cellular machinery, including reactions with DNA as well as processes like apoptosis, oxidative stress and mitochondrial disturbance. Cell viability can be evaluated either individually or concurrently with the target bioassay endpoint. For instance, one approach involves evaluating the cell viability by determining the proportion of live cells within a population. This is achieved using a fluorescent dye that selectively enters living cells, as it cannot permeate the membranes of deceased cells, resulting in fluorescence intensity directly reflecting cell viability.

As the concentration of the toxic substance approaches levels that induce cell death, the signal associated with the presumably specific toxicity of a target assay endpoint may become increasingly mixed with signals stemming from non-specific cytotoxicity burst responses [18]. Only from the observed responses of the target assay endpoint it can not be deduced what are the specific and non-specific shares in the measured signal.

Referred to as *false positive* hitcalls, these are associated with compounds where the activity response surpasses the efficacy cutoff mainly because of non-specific toxicity. Nevertheless, in many research contexts, there exists a particular interest in pinpointing specific toxicity [19]. This becomes crucial for identifying the molecular initiating event and understanding the adverse outcome pathway. Solely based on the observed signal, the challenge arises in differentiating true positives, where the compound exhibits specific toxicity without cytotoxicity interference, from false positive hitcalls. This

2.2. Chemical Target Toxicity vs. Cytotoxicity



Figure 2.3: Example of a bioassay response with cytotoxicity interference. The dotted line shows the theoretical specific toxicity effect but due to non-specific cytotoxicity (black line is cell viability), the measured effect may have an inverted U-shape within the tested concentration range. The measured specific effect may also be influenced by the presence of the cytotoxicity burst phenomenon, which can lead to a non-specific exponential growth phase before the subsequent decline in the effect curve. Figure 7.8 from [18].

introduces significant uncertainty in the reported activity hitcalls.

Nonetheless, the ToxCast pipeline is deliberately structured to minimize the occurrence of *false negative* hitcalls. The original pipeline employs a fairly inclusive risk assessment approach, ensuring that compounds with ambiguous toxicity potential are more likely to be rated as active rather than inactive. Moreover, the toxicity assessment process within this pipeline lacks proper mechanisms to differentiate between activity arising from specific and non-specific chemical toxicity.

Although not the central emphasis of this study, we investigate the possibility of reducing potential overestimation of positive hitcalls attributed to suspected non-specific components in the reported activity. This is achieved by comparing potency concentrations between the target assay endpoints and the corresponding viability or burst assay endpoints, which quantify cytotoxic cell loss or cell stress, respectively. If the probabilities indicate that a crucial potency concentration from the cytotoxicity assay endpoint is lower than that of the target assay endpoint, previously identified false positive hitcalls can be reduced by a factor reflecting the potential impact of cytotoxicity interference.

Chapter 3

Related work

The recent developments in machine learning for predicting toxicity endpoints were outlined in [20], with an observation that these advancements are primarily driven by the progress made in the field of drug discovery. The study underscores that machine learning approaches demonstrate varying performance levels across diverse toxicity endpoints, with commonly studied ones including cardiotoxicity, mutagenicity, hepatotoxicity and acute oral toxicity, but also those endpoints from the popular Tox21 data challenge [21]. The ability to predict toxicity depends significantly on the characteristics of the datasets, including differences in complexity, class distribution, and the chemical space they encompass, making it challenging to directly compare algorithm performance.

A recent study [22] explores the coverage of large-scale datasets used in machine learning for biomolecular structures, revealing their limitations in representing the full range of known structures. As the chemical space is vast, it is questionable whether the toxicity training data is an informative subset to the true distribution aimed to learn, directly challenging the fundamental assumption in machine learning. The study underscores the importance of taking into account the coverage of chemical space when assessing the effectiveness of machine learning models. In this thesis, the coverage of the chemical space was not specifically assessed, as the focus was on the performance of the models and their ability to generalize to unseen data. We recognize that the prediction ability of the developed pipeline relies on the similarity of unknown compounds to the training chemical space, but addressing this issue in detail is beyond the scope of this thesis and will be considered in future work.

Similar to MLinvitroTox, MS2Tox [23] represents another machine learning approach within the realm of predicting ecotoxicological hazards for unidentified compounds through nontarget HRMS/MS analysis. Both approaches adopt a common strategy of building their ML models based on molecular fingerprints derived from chemical structure, used to make predictions on environmental samples, utilizing fingerprints from fragmentation spectra calculated by SIRIUS+CSI:FingerID. However, ML2Tox diverges in terms of the toxicity data employed for training and testing, with its focus on toxicity

data concerning *in vivo* fish lethal concentrations from CompTox [24]. This is in contrast to MLinvitroTox, which relies on *in vitro* toxicity data from ToxCast/Tox21. Additionally, unlike MLinvitroTox, which exclusively relies on molecular fingerprints and does not utilize other physicochemical properties, MS2Tox incorporates the molecular mass of the compound as an additional feature.

In a systematic investigation using Tox21 data [25], the impact of various modeling approaches on predictive toxicology were explored, with a focus on model performance and explainability trade-offs. The study found that endpoints with higher predictability, characterized by lower data imbalance and larger datasets, performed well regardless of the modeling approach or molecular representation. For less predictable endpoints, simpler models like linear models performed similarly to complex ones, thereby emphasizing the importance of balancing predictability and interpretability. Moreover this study suggests consensus modeling and multi-task learning to enhance predictability and model performance across endpoints. In this thesis, simpler models should not be disregarded due to their higher interpretability and comparable performance. As per the recommendation, no additional investigations were pursued concerning the various molecular representations. Instead, we employed a compilation of molecular fingerprints encompassed by SIRIUS, which are already integrated into the EXPECTmine pipeline, making them inherently suitable for our purpose. Subsequently, we applied feature selection to reduce the number of relevant features. Furthermore, we implemented a consensus modeling approach, where the ultimate predictions are obtained by averaging the predictions across assay endpoints that share common attributes, such as mechanistic and biological targets.

Chapter 4

Material and Methods

4.1 Toxicity Data and Processing

4.1.1 ToxCast *invitroDB v4.1*

The most recent release of the ToxCast's database, referred to as *invitroDBv4.1*, is an extensive collection of HTS toxicity data (~100 GB). This database holds data for 10'196 compounds, each selectively tested across 1'485 assay endpoints. It employs a one-to-many relationship, where a single assay interfaces with numerous endpoints to thoroughly explore a compound's biological effects. An assay constitutes an experimental protocol for high-throughput evaluation of target molecule interactions with compounds, encompassing elements like reagents and response-measuring instruments. Meanwhile, assay endpoints involve specific measurements taken under various conditions and timeframes, delivering a comprehensive understanding of compound-induced biological responses. To get a grasp of the assay annotation structure, consult Figure 4.1.

The assays utilize a range of technologies to assess the impact of chemical compounds on a wide array of biological targets, including individual proteins, nuclear receptor signaling, developmental processes and cellular processes such as mitochondrial health. This resource originates from the collaboration of two prominent institutions: the U.S. EPA through its ToxCast program and the National Institutes of Health (NIH) via the Tox21 initiative. Using data collected from multiple research labs (refer to Table A.1 in the Appendix), this relational database is accessible to the public and can be downloaded¹ by visiting the official ToxCast website.

4.1.2 tcpl v3.0

The primary ToxCast pipeline is effectively managed through the extensive toolkit offered by the `tcpl`² package, which includes a variety of tools for high-throughput

¹<https://www.epa.gov/chemical-research/exploring-toxcast-data>, released on Sept 21, 2023

²<https://github.com/USEPA/CompTox-ToxCast-tcpl>

4.1. Toxicity Data and Processing

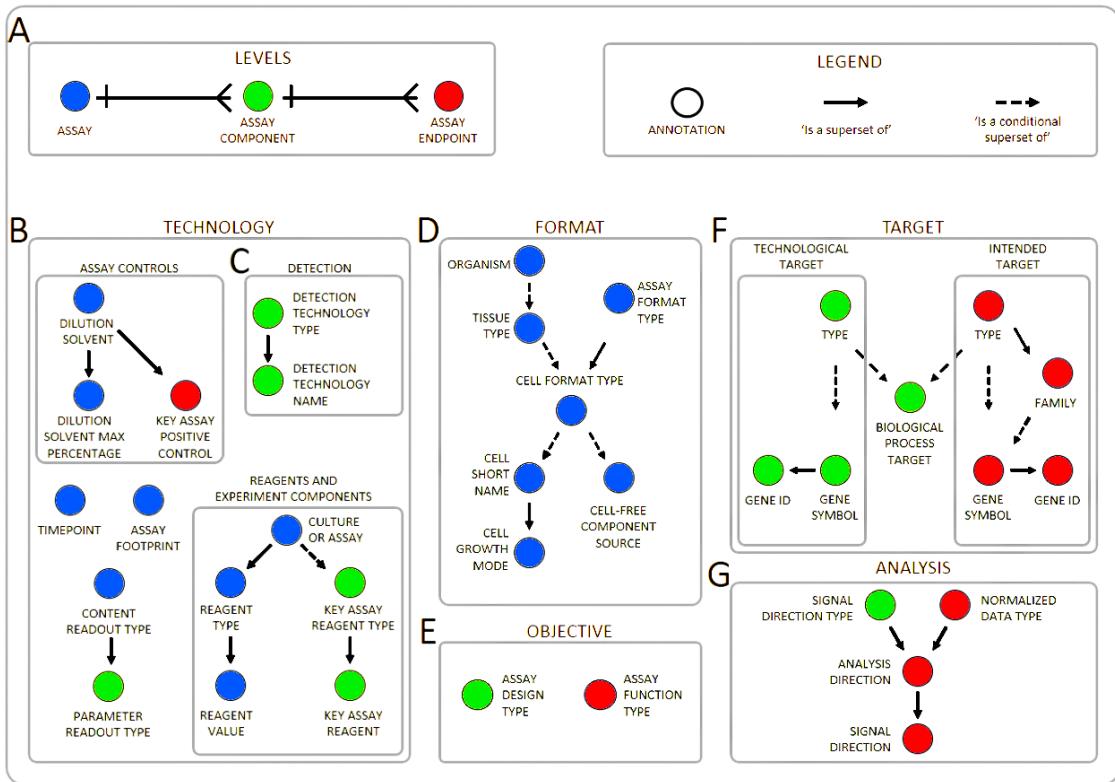


Figure 4.1: The annotations for assay endpoints include (A) information for identifying the (color-coded) assay entity, (B) design-related data, (C) details about the target, and (D) information regarding the analysis. These annotations exhibit relationships that can be either one-to-many or conditional, with some dependencies not being applicable in certain cases. This modified figure is sourced from [26].

screening data management. It enables reproducible concentration-response modeling and populates the MySQL database, invitroDBv4.1. The multiple concentration screening paradigm intends to pinpoint the bioactivity of compounds, while also estimating their efficacy and potency. In Section 4.2, we introduce `pytcp1`, a Python reimplementation of the major components that underpin the entire ToxCast pipeline. It should be noted that these components, as presented in the following, are applicable to both `tcp1` and `pytcp1`.

4.1.3 Concentration-Response Series

Each compound c_j tested within an assay endpoint a_i involves the collection of the respective *concentration-response series* (CRS) denoted as CRS_{ij} , showcased in Figure 4.2. A CRS is represented as a set of concentration-response pairs:

$$CRS_{ij} = \{(conc_{1_{ij}}, resp_{1_{ij}}), (conc_{2_{ij}}, resp_{2_{ij}}), \dots, (conc_{n_{\text{datapoints}_{ij}}}, resp_{n_{\text{datapoints}_{ij}}})\}$$

where $n_{\text{datapoints}_{ij}}$ varies based on the number of concentrations tested.

4.1. Toxicity Data and Processing

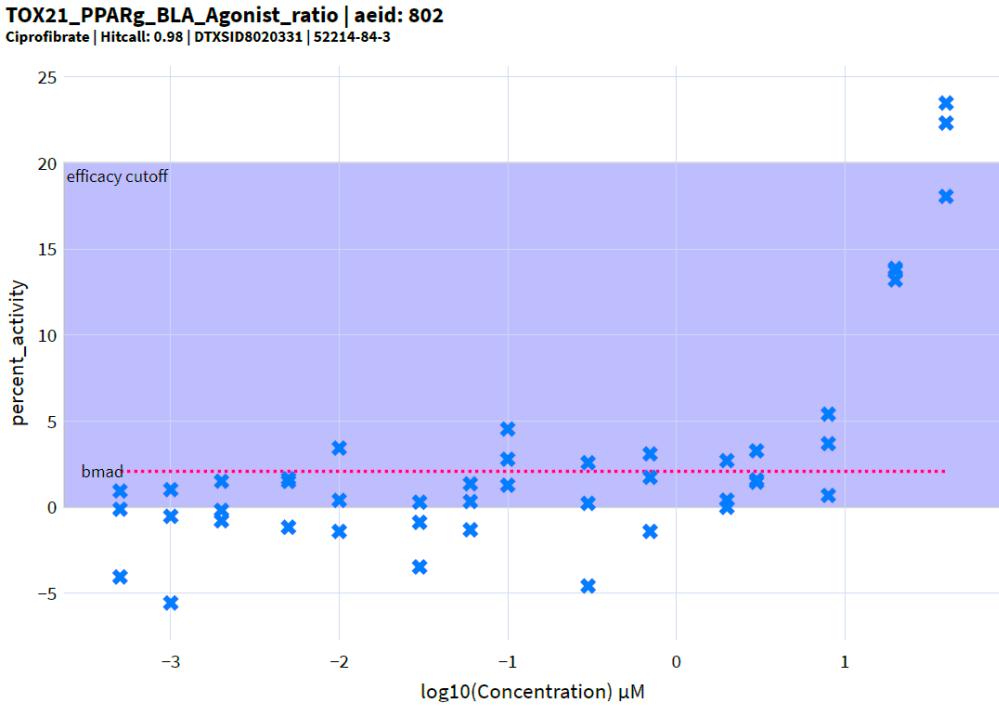


Figure 4.2: The CRS belongs to *Ciprofibrate* (DTXSID8020331), tested in the assay endpoint TOX21_PPARg_BLA_Agonist_ratio (aeid=802). The shaded region represents the estimated efficacy cutoff, and the dashed line represents the baseline median absolute deviation (BMAD), as explained in Section 4.1.4.. This particular series comprises a total of $k = 45$ concentration-response pairs and is structured into $n_{conc} = 15$ distinct concentration groups, with each group consisting of $n_{rep} = 3$ replicates.

Table 4.1: Description of Parameters

Quantity	Description
$n_{datapoints_{i,j}}$	Total number of concentration-response pairs ($ CRS $)
$n_{groups_{i,j}}$	Number of distinct concentrations tested
$n_{replicates_{i,j}}$	Number of replicates for each concentration group
$min_{conc_{i,j}}$	Lowest concentration tested
$max_{conc_{i,j}}$	Highest concentration tested

In practice, concentrations are often subjected to multiple testing iterations, resulting in the distinct concentration groups with replicates. Table 4.1 presents the key quantities associated with an individual CRS when considering a specific assay endpoint a_i and compound c_j . To visualize the variations in these quantities across the complete set of analyzed CRS in this work, please refer to Figure A.1 in the Appendix.

Concentration-response pairs, along with essential sample information such as well type and assay well-plate indices, can be retrieved by combining tables $mc0$, $mc1$, and

4.1. Toxicity Data and Processing

mc3 from invitroDBv4.1, which represent the raw data. A special role is assigned to the control wells, which typically contain untreated samples or samples with a known, non-toxic response. They are used as a baseline to normalize the treated samples and account for any background noise in the assay [27]. The concentrations are transformed to the logarithmic scale in micromolar, while the responses are control well-normalized to either fold-induction or percent-of-control activity:

1. **Fold Induction:** is a measure used to quantify how much, for instance, gene expression has changed in response to a treatment compared to its baseline level from the control well set. E.g., if a gene is expressed five times higher in a treated sample compared to the control, the fold induction would be 5.
2. **Percent of Control:** is another way to express the relative change in an activity due to a treatment compared to the control.

4.1.4 Efficacy Cutoff

The evaluation of a compound's bioactivity is significantly influenced by the specific *efficacy cutoff* associated with each assay endpoint, as exemplified for a tested compound in Figure 4.2. It serves as a threshold that differentiates active and inactive compounds, essentially defining the minimum response level of toxicity that is biologically relevant. The process of establishing this threshold involves estimating the noise level in the assay endpoint based on the baseline median absolute deviation (BMAD). The BMAD is calculated using baseline response values, which are assumed to come from either untreated control wells or test samples from the two lowest concentrations. This calculation is performed just once for the entire assay endpoint.

4.1.5 tcplFit2

*tcplFit2*³ is an extension to *tcpl*, focused on curve-fitting and hit-calling. The package also offers a flexible and robust fitting procedure, allowing for the use of different optimization algorithms and the incorporation of user-defined constraints. This sets it apart from other open-source CRS modeling packages such as *drc* and *mixtox*, as it is explicitly designed for HTS concentration-response data.

4.1.6 Curve Fitting

Curve fitting involves the adjustment of mathematical models to best match observed data. Maximum Likelihood Estimation (MLE) is a statistical technique used to find parameter values that maximize the likelihood of observing the data within the model. The likelihood function represents the probability of observing the data given the model's parameters. The various curve fitting models investigated in *tcplFit2* are summarized in Table 4.2 and visually presented in Figure 4.3.

³<https://github.com/USEPA/CompTox-ToxCast-tcplFit2>

4.1. Toxicity Data and Processing

Table 4.2: tcplFit2 Model Details

Model	Label	Equations ¹
Constant	constant	$f(x) = 0$
Linear	poly1	$f(x) = ax$
Quadratic	poly2	$f(x) = a \left(\frac{x}{b} + \left(\frac{x}{b} \right)^2 \right)$
Power	power	$f(x) = ax^p$
Hill	hill	$f(x) = \frac{tp}{1 + \left(\frac{ga}{x} \right)^p}$
Gain-Loss	gnls	$f(x) = \frac{tp}{\left(1 + \left(\frac{ga}{x} \right)^p \right) \left(1 + \left(\frac{x}{la} \right)^q \right)}$
Exponential 2	exp2	$f(x) = a \left(\exp \left(\frac{x}{b} \right) - 1 \right)$
Exponential 3	exp3	$f(x) = a \left(\exp \left(\left(\frac{x}{b} \right)^p \right) - 1 \right)$
Exponential 4	exp4	$f(x) = tp \left(1 - 2^{-\frac{x}{ga}} \right)$
Exponential 5	exp5	$f(x) = tp \left(1 - 2^{-\left(\frac{x}{ga} \right)^p} \right)$

¹ Constrained parameters: a : x-scale, b : y-scale p : gain power, q : loss power, tp : top, ga : gain (x), la : loss (x)



Figure 4.3: Figure obtained from [28].

4.1. Toxicity Data and Processing

In all models, it is presumed that the errors conform to a Student's t -distribution [28]. The presence of heavier tails in this t -distribution reduces the impact of outlier values, resulting in more resilient estimates compared to the frequently employed normal distribution. This robust model fitting approach eliminates the requirement for filtering out potential outliers before the fitting process. For a comprehensive explanation of the fitting procedure, please consult the official vignette⁴.

The *Akaike Information Criterion (AIC)* serves as the metric for assessing the goodness of fit of models, defined by the formula: $AIC = -2 \log(L(\hat{\theta}, y)) + 2K$, where $L(\hat{\theta}, y)$ is the likelihood of the model θ given the data and K is the number of model parameters. The model with the lowest AIC value is chosen as the *winning* model. The winning model is then used to estimate the efficacy and potency of the compound. The potency estimates, also called *point-of-departure (POD)* estimates, are derived from the fitted curve, identifying certain *activity concentrations (AC)* at which the curve first reaches certain response levels. Central POD estimates are depicted graphically in Figure 4.4.

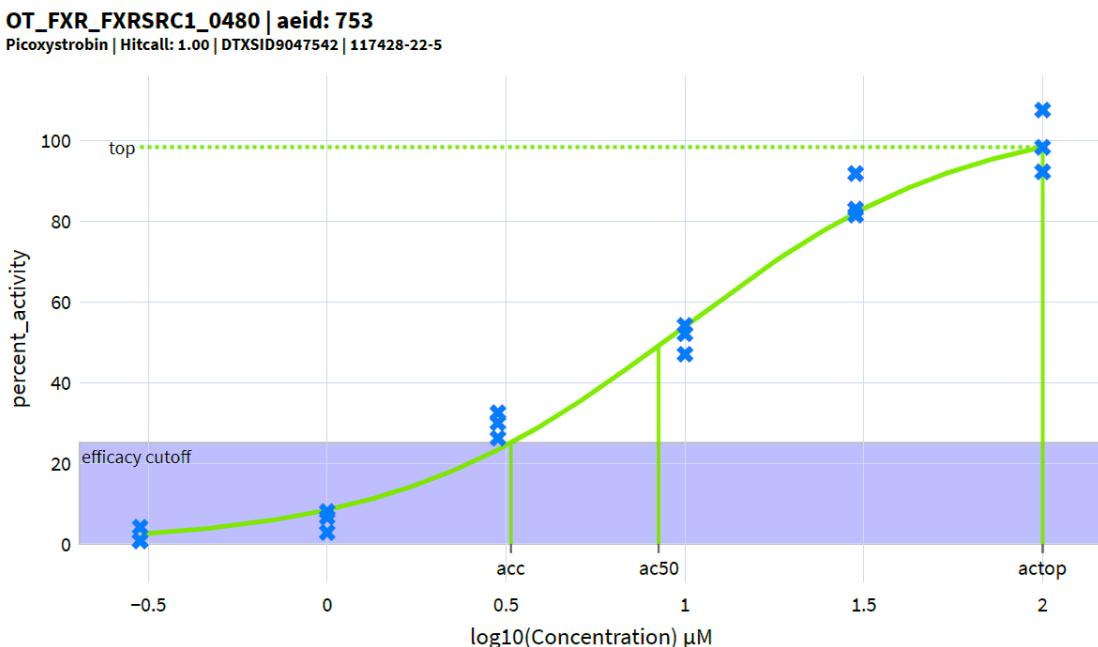


Figure 4.4: The Point of Departure (POD) potency estimates for the chemical compound *Picoxystrobin* (DTXSID9047542) tested in the assay endpoint with *aeid* = 753. The efficacy cutoff is defined at ~ 25 percent-of-control activity. The winning fit model was the Hill function. ACC: The AC at the efficacy cutoff is $3.3 \mu M$. AC50: The AC at 50% of the maximum response is $8.4 \mu M$. ACTop: The AC at the maximum response is $100 \mu M$.

⁴https://cran.r-project.org/web/packages/tcpl/vignettes/Data_processing.html

4.1. Toxicity Data and Processing

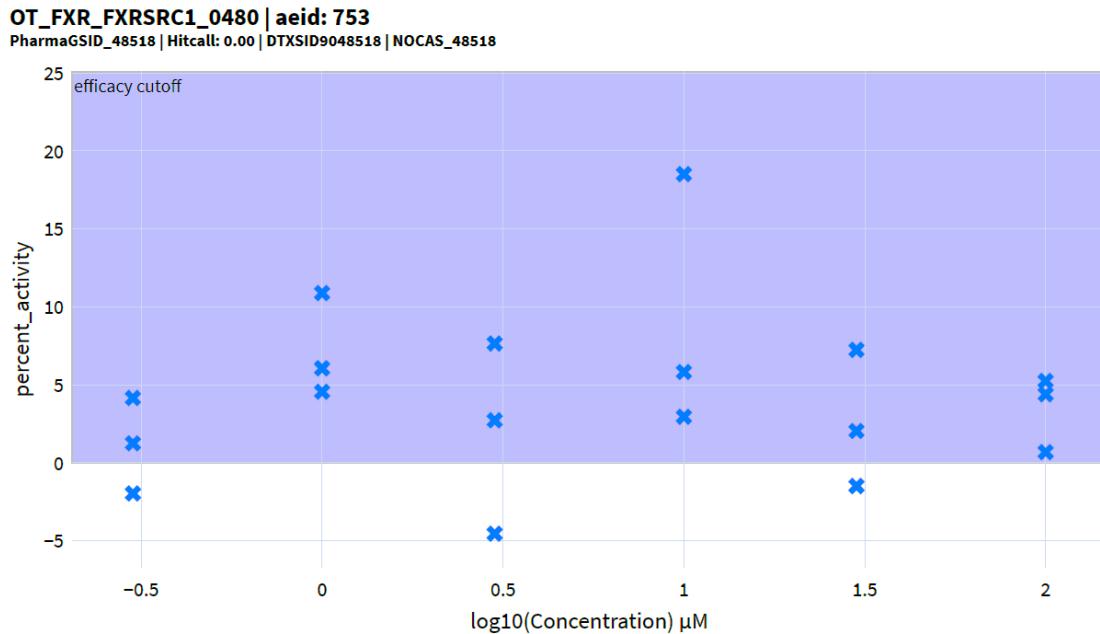


Figure 4.5: The Point of Departure (POD) potency estimates are not available for the chemical compound *PharmaGSID_48518* (DTXSID9048518), also tested in the assay endpoint with *aeid* = 753. In this situation, it was unnecessary to fit curves as no response reached or exceeded 80% of the efficacy cutoff, clearly indicating the inactivity of the compound. In such scenarios, a calculation of POD estimates is not applicable.

4.1.7 Hit Calling

A *binary hitcall* hitcall simplifies the classification of the estimated activity for a specific compound tested in a particular assay endpoint and results in toxicity testing in a *toxic* (hit) or *non-toxic* (no hit) classification.

The *continuous hitcall*, on the other hand, provides a more nuanced evaluation of the likelihood that a compound is active. The *tcplFit2* package introduces a continuous hitcall, based on the product of the following three distinct probability values [27]:

- i. that at least one median response is greater than the efficacy cutoff, computed by using the error parameter from the model fit and Student's *t*-distribution to calculate the odds of at least one response exceeding the efficacy cutoff;
- ii. that the top of the winning fitted curve is above the cutoff which is the likelihood ratio of the one-sided probability of the efficacy cutoff being exceeded;
- iii. that the winning AIC value is less than that of the constant model:

$$\frac{e^{-\frac{1}{2}AIC_{winning}}}{e^{-\frac{1}{2}AIC_{winning}} + e^{-\frac{1}{2}AIC_{cnst}}} \quad (4.1)$$

In certain instances, compounds underwent multiple tests within a single assay endpoint, leading to their association with multiple CRS. In these exceptional cases, a hitcall is

4.2. New Toxicity Pipeline Implementation: pytcpl

computed for each CRS, and then the highest hitcall value is recorded as the compound's ultimate hitcall.

4.1.8 Flagging

Finally, after processing, each CRS is assigned to an appropriate fit category based on the level of certainty in the estimated bioactivity. Additionally, cautionary flags are assigned to account for problematic data series or uncertainty related fits and hits.

4.2 New Toxicity Pipeline Implementation: pytcpl

4.2.1 Introduction

This thesis introduces `pytcpl`⁵, a streamlined Python repository inspired by the R packages `tcp1` and `tcp1Fit2`. This package was developed to accomodate customizable processing steps and facilitate interactive data visualization with an own *Curve Surfer*⁶. The package optimizes data storage and generates compressed Parquet files of the relevant raw data and metadata from *invitroDBv4.1*. Exclusively utilizing this repository eliminates the need for a complex and extensive database installation, rendering downstream analysis more accessible and efficient. It enables researchers who prefer Python to easily participate in data analysis and exploration, overcoming limitations associated with using R code.

The `pytcpl` pipeline adds a setup and post-processing step around the main pipeline:

- **Setup:** This step involves user-specified subsetting of assay endpoints, tagging assays with external assay annotations, enabling workload balancing for distributed processing and generating Parquet files from all raw and metadata, optionally for database decoupled analysis.
- **Main** (similar to `tcp1+tcp1Fit2`): This step involves cutoff determination, curve fitting, hit calling and flagging.
- **Post-Processing:** This step has the goal of improving the overall quality of the data and involves post-processing curation, cytotoxicity interference reevaluation and the custom export of the final results.

4.2.2 Setup step

Subsetting Data

For a better data comprehension, the presence matrix denoted as $P \in \{0,1\}^{m \times n}$ is introduced. In this matrix, rows (indexed by i) represent assay endpoints a_i , and columns (indexed by j) indicate whether testing was performed (1) or not performed (0)

⁵<https://github.com/rbBosshard/pytcpl>

⁶<https://pytcpl.streamlit.app/>

4.2. New Toxicity Pipeline Implementation: pytcp1

for compound c_j in those endpoints. Because of selective compound testing, the matrix P is sparse. For a visual representation of this presence matrix encompassing all assay endpoints and compounds in *invitroDBv4.1*, refer to Figure 4.6. Note that the presence matrix was structured by the ranking the the number of compounds associated with each assay endpoint, with compounds sorted in descending order of their occurrence frequency.

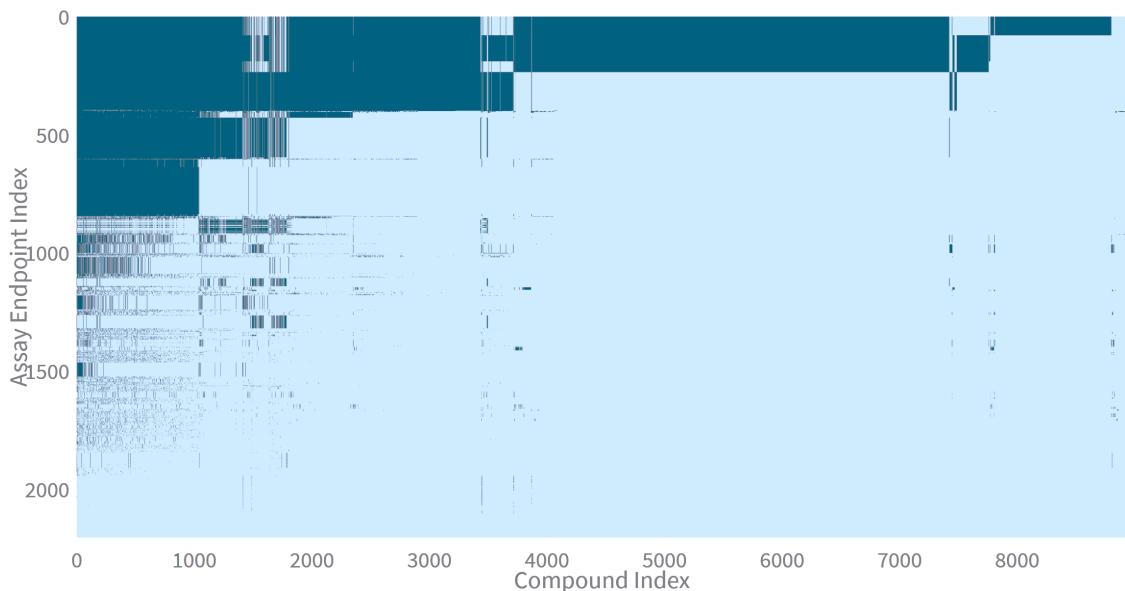


Figure 4.6: The presence matrix P_{all} , covers all assay endpoints and compounds from *invitroDBv4.1*, totaling $m = 2'205$ assay endpoints and $n = 8'935$ compounds, excluding 606 compounds lacking molecular fingerprints. There are 3'196'178 concentration-response series available.

However, we do not utilize the whole dataset and exclude *biochemical in vitro* assays and focus solely on *cell-based in vitro* assays that use intact cells (e.g. human primary cells and cell lines and rat primary liver cells) to assess cellular responses when exposed to test substances. Moreover, we exclusively considered assay endpoints that have been tested with a minimum of 1000. This selection criterion ensures the presence of adequate data for subsequent training of robust machine learning models. You can refer to Figure 4.7 for a visual representation of the presence matrix P which includes only this particular subset of assay endpoints. From this moment forward, we will refer to this specific subset as *the dataset*, which will be the focus of this thesis.

External Assay Annotation

The investigated assay endpoints are enriched with external annotations attributed by the Integrated Chemical Environment (ICE) [29], which provide valuable context and information about each endpoint. These annotations encompass the following aspects:

4.2. New Toxicity Pipeline Implementation: pytcp1

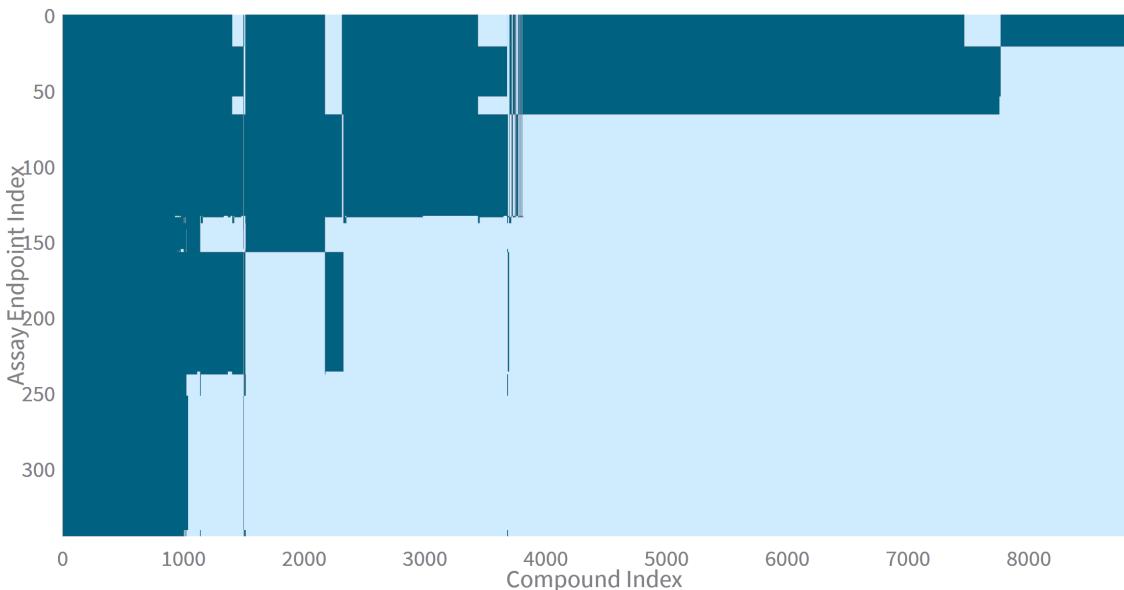


Figure 4.7: P_{subset} covers a specific subset of relevant assay endpoints and compounds considered for this thesis, totaling $m = 345$ assay endpoints and $n = 8'804$ compounds. Assay endpoints with less than 1'000 compounds tested were omitted. There are 1'043'222 concentration-response series available.

1. **Toxicity Endpoint:** This annotation specifies the type of toxicity or adverse effect associated with each assay endpoint, helping to clarify the specific aspect of toxicity under investigation.
2. **Mechanistic Target:** This annotation sheds light on the particular target mechanism or biological pathway being studied.
3. **Mode of Action:** The annotations also describe how the tested compounds interact with the mechanistic targets and provides insights into the underlying biological processes or actions involved.

4.2.3 Main step

The pytcp1 main pipeline is similar to the R-based `tcp1+tcp1Fit2` pipeline, with the exception of the curve fitting stage where the pytcp1 pipeline made a notable modification by including a novel model the removal of the Exponential 3 model. The exclusion is a consequence of its infrequent selection as the winning model within the `tcp1` pipeline, indicating limited effectiveness in model fitting, as outlined in [30]. Furthermore, an additional Gain-Loss 2 model was introduced during the curve-fitting stage. This model has one model parameter less than the Gain-Loss 1 model, which helps mitigate the risk of overfitting CRS data, making it less susceptible to outliers. Table 4.3 provides an overview of these changes within the pytcp1 pipeline for reference.

4.2. New Toxicity Pipeline Implementation: pytcp1

Table 4.3: pytcp1 Model Updates

Model	Label	Equations ¹	Role in pytcp1
Exponential 3	exp3	$f(x) = a \left(\exp \left(\left(\frac{x}{b} \right)^p \right) - 1 \right)$	Omitted
Gain-Loss 2	gnls2	$f(x) = \frac{tp}{1 + \left(\frac{ga}{x} \right)^p} \exp(-qx)$	New

4.2.4 Post-Processing step

Post-Processing Curation

Following the ICE guidelines⁷, quality filters were implemented to enhance the processed concentration-response series. This step introduces OMIT/PASS warning flags, which could be applied based on assay endpoints or compound quality control criteria.

Cytotoxicity Interference Reevaluation

As previously discussed in Chapter 2, the assessment of compound toxicity can be complicated by the presence of non-specific cytotoxic responses. In this section, we delve into the exploration of a method for reevaluating the reported hitcall status of active compounds, considering the estimated extent of cytotoxicity interference. The cytotoxicity of a compound in a target assay endpoint may be assessed by comparing the activity concentration at the efficacy cutoff, represented as ACC_{target} , with that of its corresponding viability assay endpoint counterpart (as exemplified in Table 4.4), referred to as ACC_{cyto} . If no counterpart is available in the database, we apply a statistical approach to calculate a cytotoxicity estimate. It uses the median ACC for the compound of interest across a set of assay endpoints dedicated for capturing the cytotoxicity burst. For both cases, the ACC is assumed to have a Gaussian error distribution. Cytotoxicity in terms of the respective potencies is assumed when: $ACC_{cyto} \leq ACC_{target}$. Thus, the probability of a compound being cytotoxic can be determined by:

$$P(\text{cytotoxic}) = P(ACC_{cyto} - ACC_{target} \leq 0) = \Phi \left(\frac{ACC_{cyto} - ACC_{target}}{\sqrt{SD_{ACC_{cyto}}^2 + SD_{ACC_{target}}^2}} \right)$$

where Φ is the Gaussian cumulative distribution function. The standard deviations $SD_{ACC_{cyto}}$ and $SD_{ACC_{target}}$ are unknown but are estimated as $0.3 \log_{10} \mu M$ units [31]. For the statistical approach with the cytotoxicity burst assays, $SD_{ACC_{cyto}}$ can be derived from the median absolute deviation (MAD) of the respective ACC values. Additionally, $P(\text{cytotoxic})$ is multiplied with the ratio of the number of cytotoxicity burst assay endpoints in which the compound exhibited activity (n_{hit}) to the total number of cytotoxicity burst assay endpoints in which the compound was tested (n_{tested}). Ultimately, $P(\text{cytotoxic})$ is multiplied with the original continuous hitcall of

⁷<https://ice.ntp.niehs.nih.gov/DATASETDESCRIPTION?section=cHTS>

4.2. New Toxicity Pipeline Implementation: pytcpl

active compounds. The final cytotoxicity-corrected hitcall is then defined as follows:
 $\text{hitcall}_{\text{cyto-corrected}} = \text{hitcall}_{\text{original}} * (1 - P(\text{cytotoxic}))$.

Table 4.4: Each assay endpoint has an assay identifier (aid) used to match it with its viability counterpart that assesses cell loss. In this example, *APR_HepG2_CellLoss_24hr* (aid=26) matches with aid=38 and aid=40. Similarly, *APR_HepG2_CellLoss_72hr* (aid=46) matches with aid=58 and aid=60.

aeid	assay endpoint name	aid	assay name	assay function type
26	APR_HepG2_CellLoss_24hr	3	APR_HepG2_24hr	viability
38	APR_HepG2_P-H2AX_24hr	3	APR_HepG2_24hr	signaling
40	APR_HepG2_p53Act_24hr	3	APR_HepG2_24hr	signaling
46	APR_HepG2_CellLoss_72hr	4	APR_HepG2_72hr	viability
58	APR_HepG2_P-H2AX_72hr	4	APR_HepG2_72hr	signaling
60	APR_HepG2_p53Act_72hr	4	APR_HepG2_72hr	signaling

4.2.5 Curve Surfer

Figure 4.8 presents the developed *Curve Surfer*, a browser-based application that enables interactive data exploration and visualization of the processed data. The curve surfer tool is built using Streamlit, an open-source Python library that makes it easy to build custom web-apps for machine learning and data science.

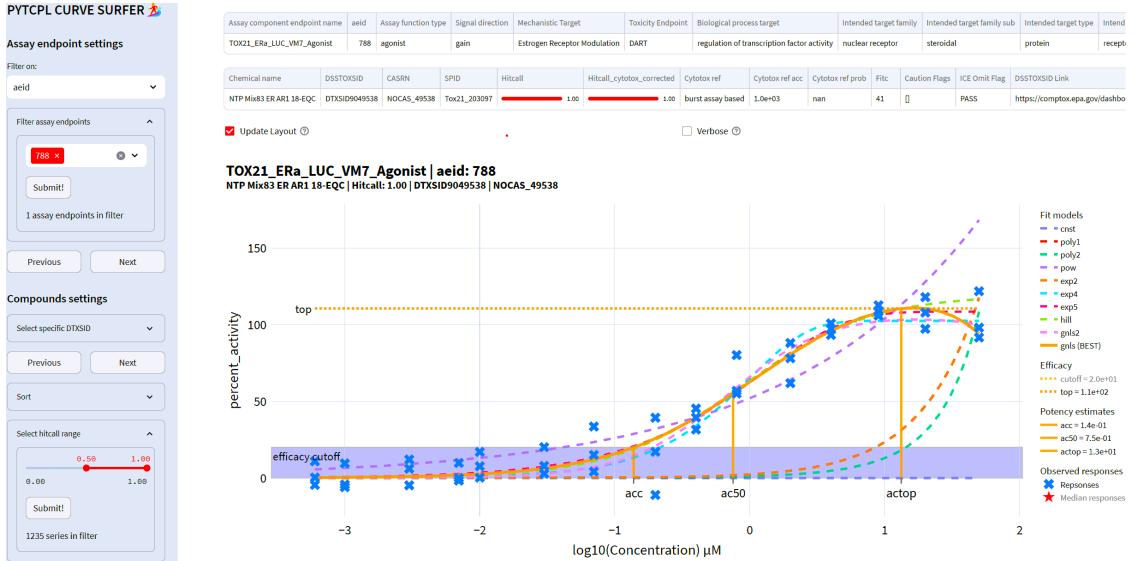


Figure 4.8: The curve surfer provides the capability to narrow down assay endpoints based on critical annotations. Compounds can be selectively filtered using their DTXSID. Users can navigate through assay endpoints or compounds within the current assay endpoint but can also be filtered by their hitcall value or POD estimates using a range slider. The tool presents the CRS data, curve fit models and associated metadata.

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

4.3 Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

We employed distinct machine learning models for each assay endpoint, enabling the prediction of compound toxicity unique to each assay endpoint, coleveraging molecular structure inputs. These models utilize molecular fingerprints to predict compound toxicity, as illustrated in Figure 4.9. We applied binary classification and regression models. In the case of binary classification, we use a threshold of 0.5 to convert the continuous hitcall target values into binarized outcomes.

To create these individual datasets, we extract compounds that possess toxicity data for a given assay endpoint from the outputs generated by `pytcp1`. The associated hitcall values for these compounds are used as target variables within the machine learning model.

The binary input features for the model are molecular fingerprints with 2362 bits derived from chemical structures. The structural data was obtained from the U.S. EPA's DSSTox database, accessed through the CompTox Chemicals Dashboard⁸. The structural data mining and the necessary structure cleanup, as mentioned in Section 2.1, was conducted by Dr. Kasia Arturi⁹.

Index	X (fingerprint features)						y (activity label)
Compound ID	fps 1	fps 2	fps 3	fps 4	...	fps n	hitcall
DTXSID 1	0	1	0	0	...	0	0.00
DTXSID 2	1	0	0	0	...	1	0.01
DTXSID 3	0	0	1	1	...	0	0.00
DTXSID 4	1	0	0	0	...	0	0.98
...
DTXSID m	0	0	1	0	...	1	0.01

Figure 4.9: Schematic example of a machine learning dataset related to a single assay endpoint. The dataset is structured into a feature matrix with $n = 2362$ and a target vector. The feature matrix consists of molecular fingerprints, and the target vector is the hitcall value. For binary classification, the hitcall value is binarized based on a specific activity threshold ($=0.5$)

The machine learning pipeline is structured into three main stages: model training, model evaluation and model application. The following sections and the diagram illustrated in Figure 4.10 provide a detailed description of each stage.

⁸<https://www.epa.gov/comptox-tools>

⁹<https://gitlab.renkulab.io/expectmine/generating-fingerprints>

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

MLinvitroTox

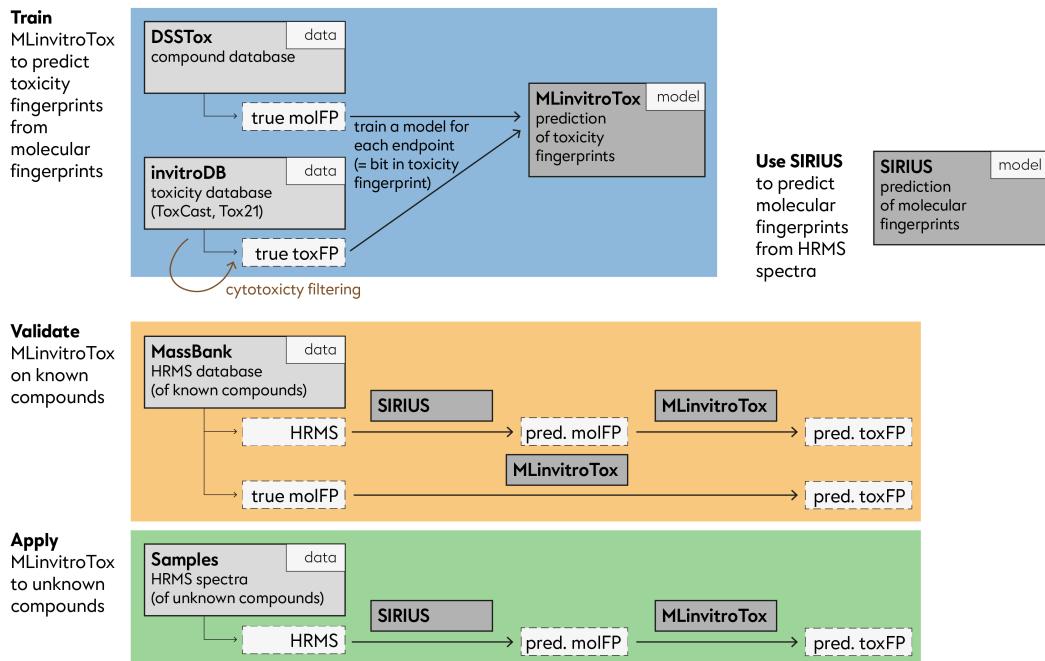


Figure 4.10: MLinvitroTox: Machine Learning Pipeline Steps. Figure created by Lili Gasser.

4.3.1 Training

The train stage is summarized in Figure 4.11 and involves the generation of individual machine learning models for each assay endpoint. Each model is trained on a subset of the dataset with an 80/20 train-validation split.

Feature Selection

We applied feature selection¹⁰ using either a xgboost or random forest model to narrow down the assay-endpoint relevant fingerprint features. The number of selected features is based on those features that exceed the mean feature importance threshold. All subsequent estimators are trained on these selected features.

Model Selection

The following supervised machine learning models from the `sklearn` library were considered:

¹⁰`sklearn.feature_selection.SelectFromModel`

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

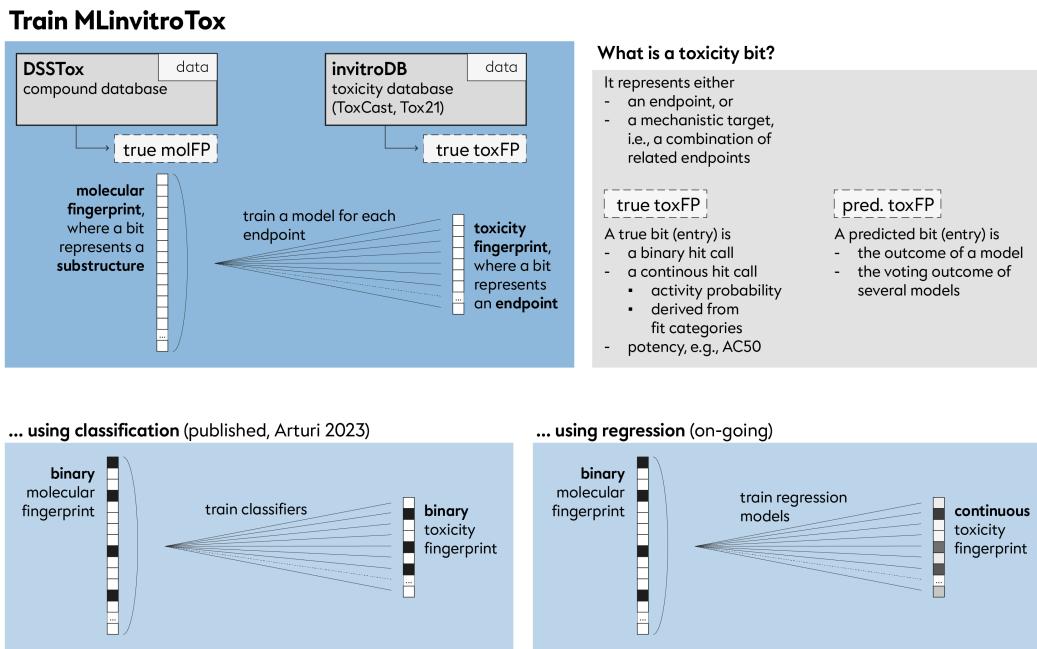


Figure 4.11: MLinvitroTox Train Step. Figure created by Lili Gasser.

1. **Logistic Regression** is a linear model that utilizes the logistic function to model binary dependent variables. It serves as a straightforward and interpretable model, often employed as a baseline for binary classification tasks.
2. **Support Vector Machine** is a robust model with a lower susceptibility to overfitting and the ability to handle high-dimensional feature spaces.
3. **Random Forest** is a bagging (bootstrap aggregating) ensemble learning technique that constructs a multitude of decision trees during training and combines their predictions, resulting in robust and accurate models with the advantage of reduced overfitting and the ability to handle high-dimensional data.
4. **XGBoost** is a gradient boosting ensemble learning technique that combines multiple weak learner decision trees sequentially, with each new learner giving more weight to the examples that the previous learners struggled with. It provides typically high predictive accuracy and efficiency through techniques like gradient optimization and regularization.
5. **Multi-Layer Perceptron** is a type of artificial neural network that consists of multiple layers of interconnected neurons and is used for various machine learning

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

tasks, offering the advantage of modeling complex non-linear relationships in data.

For every machine learning model, the selection process is based on a grid search over a set of hyperparameters, using 5-fold cross-validation. The hyperparameters, specified in a separate config file, are optimized for binary classification based on the F_β score, a generalization of the F_1 . The F_1 -score is the harmonic mean of the precision and recall and the more generic F_β score applies additional weights, valuing one of precision or recall more than the other. We set $\beta = 2$ to value recall higher than precision.

4.3.2 Evaluation

To assess the performance of our trained models, we used two separate validation sets that were not part of the training data:

1. The first validation set, referred to as *internal validation set*, was used to evaluate how well the best estimator found by the grid search 5-fold cross-validation generalizes for unseen compounds. This set was randomly drawn from the tested compounds in the particular assay endpoints, ensuring that the number of active and inactive compounds was balanced. Figure 4.12 depicts the hitcall counts (support) across the respective assay endpoints, emphasizing the prevalent imbalance.

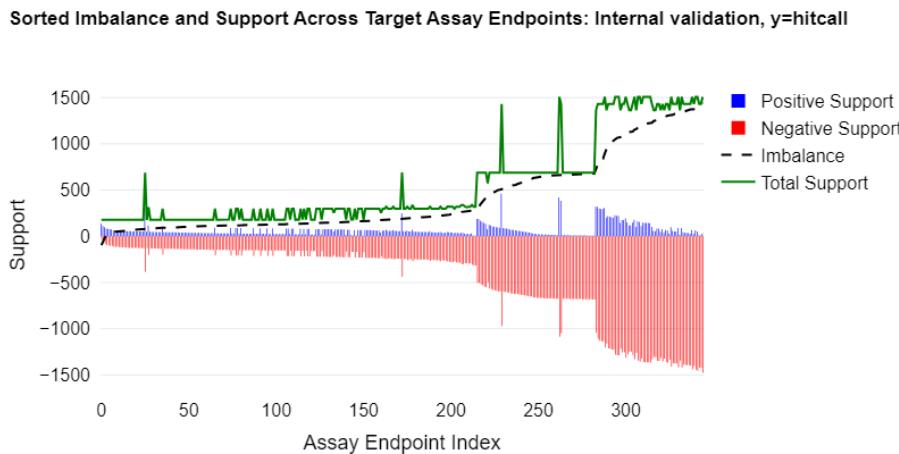


Figure 4.12: Support and imbalance in the internal validation set across the target assay endpoints.

2. The *MassBank validation set*, as the second validation set, was used to evaluate the model's generalization capabilities, focusing on the domain gap between chemical structure space and fragmentation spectra. The MassBank validation set comprises compounds with both actual (from known structure) and SIRIUS-

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

predicted fingerprints originating from MassBank spectra data. Having such a paired fingerprint dataset allows for a comparative assessment of model performance concerning actual and predicted fingerprints. However, the accuracy of SIRIUS framework's fingerprint predictions plays a critical role in assessing how well the model performs when applied to environmental samples. Figure 4.13 illustrates the respective fingerprint dissimilarities for compounds in the Massbank validation set.

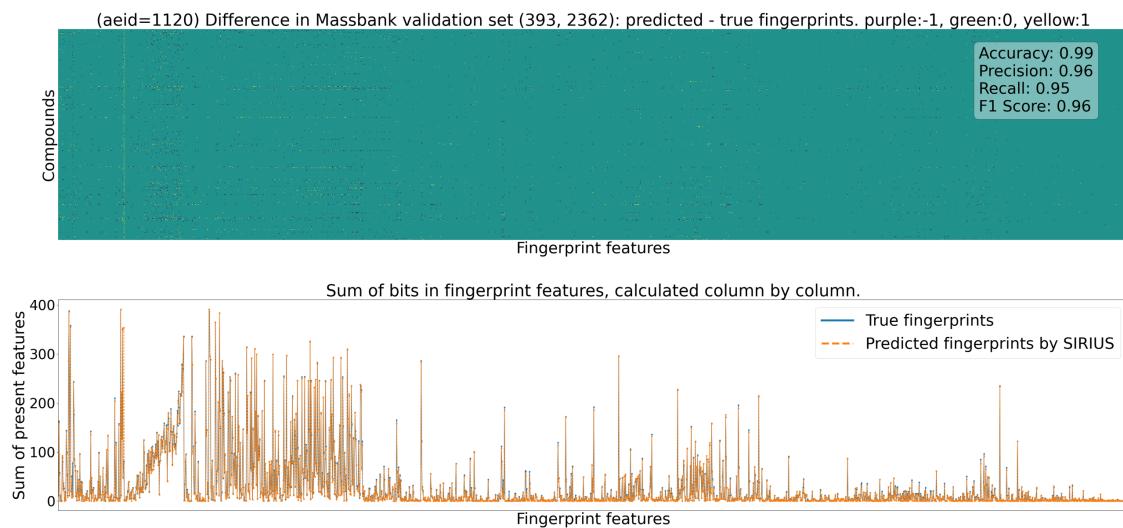


Figure 4.13: The first figure shows the molecular fingerprint dissimilarities for compounds in the Massbank validation set, illustrated for the assay endpoint TOX21_FXR_BLA_Antagonist_ratio with aeid=1120. The dissimilarity matrix displays the subtraction of SIRIUS-predicted fingerprints from the actual chemical structure fingerprints. In the second figure, it becomes apparent that predicting features that are more commonly present in the fingerprint across compounds poses a greater challenge for the SIRIUS framework.

It is important to highlight that compounds susceptible to data leakage (listed here¹¹), which were previously used in training the SIRIUS+CSI:FingerID prediction model, have been excluded. If not excluded, the model's predictive capability for a compound's fingerprint would likely be very high because the model has already seen and learned from these data points during training. In the worst-case scenario, it could have essentially memorized the training data and the evaluation would become meaningless.

After excluding these compounds, we are left with a potential upper limit of 429 compounds that can be safely used for MassBank validation, as shown in Figure 4.14. Note that the size of the MassBank validation set varies across different assay endpoints due to variations in the overlap between the compounds available in MassBank and those assessed in the assay endpoints, as depicted in Figure 4.15a.

¹¹<https://bio.informatik.uni-jena.de/software/sirius/>

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline

Note that the models have been retrained for the MassBank validation set in order to include the distinct compounds from the internal validation set. Furthermore, Figure 4.15b provides an illustration of the distribution of hitcall representation compared to the *validation* set.

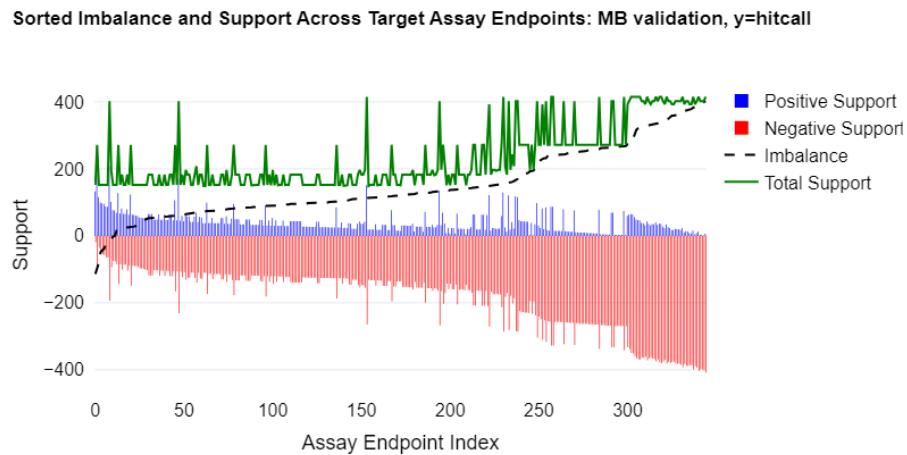


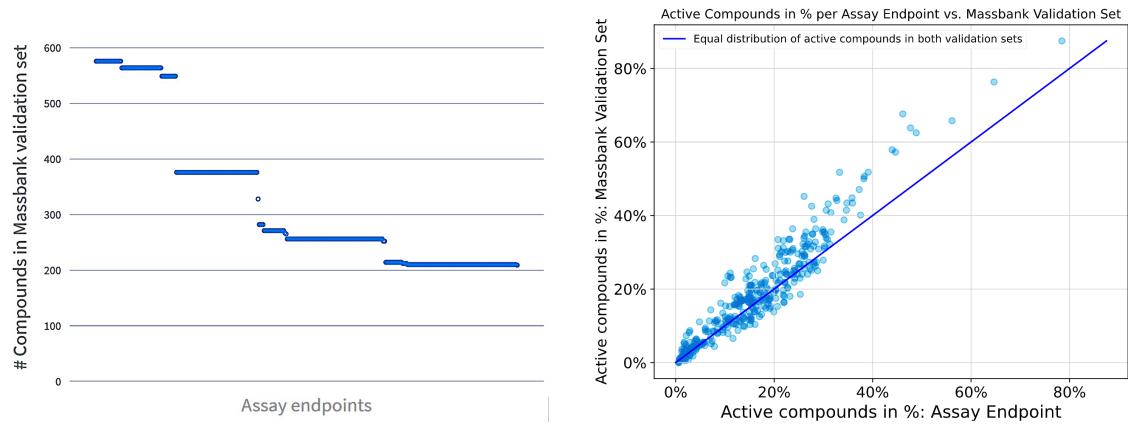
Figure 4.14: Support and imbalance in the internal validation set across the target assay endpoints.

For subsequent predictions, we retrained the models on the full dataset.

4.3.3 Application

The presence of distinct prediction models for each assay endpoint enables the grouping of these endpoints based on their annotations, such as the biological process or the mechanistic target annotation. This ultimately results in toxicity predictions averaged within these groups. These collective toxicity predictions are referred to as *toxicity fingerprint* which serve the purpose of identifying the most toxic compounds for each specific assay endpoint.

4.3. Toxicity Prediction from Molecular Fingerprints: Machine Learning Pipeline



(a) This graphic depicts the number of compounds, across various assay endpoints, for which we have both toxicity data and SIRIUS-predicted fingerprints derived from MassBank spectra. The range in the number of compounds in the MassBank validation set, which spans from 149 to 415, is due to differences in the overlap with the compounds tested in the respective assay endpoints.

(b) This figure plots the ratio of active (binarized) hitcall values for all compounds tested in the assay endpoint (x-axis) against the ratio of active (binarized) hitcall values for the compounds in the MassBank validation set (y-axis). The blue line represents the ideal case where the ratios are equal, and the validation set perfectly mirrors the entire dataset with respect to the target variable.

Figure 4.15: MassBank validation set.

Chapter 5

Results

In the scope of this thesis, the main focus was on binary classification models. While we did explore regression models, the primary findings are derived from evaluation of the classification models and thus omitted the regression results.

5.1 Binary Classification

In the context of prioritizing compounds based on hazard assessment, the objective is to maximize the probability of detecting toxic compounds while minimizing the the number of false alarms, which represent non-hazardous compounds misclassified as toxic. Initially, we introduce the performance metrics used for model evaluation.

5.1.1 Evaluation Metrics

When assessing the performance of a binary prediction model with a validation dataset containing known target values, four key quantities come into play (and presented as Confusion Matrix in Table 5.1):

Table 5.1: Confusion Matrix

	Actual Negative	Actual Positive
Predicted Negative	TN	FN
Predicted Positive	FP	TP

- True Positives (TP): The number of correctly predicted active cases.
- True Negatives (TN): The number of correctly predicted inactive cases.
- False Positives (FP): The number of incorrectly predicted active cases.
- False Negatives (FN): The number of incorrectly predicted inactive cases.

5.1. Binary Classification

The classification threshold is a crucial parameter dictating how the model assigns data points to one of two classes based on predicted class probabilities. This threshold significantly influences the model's metrics, as illustrated in Figure 5.1 and showcased in Figure 5.2 and Figure 5.3.

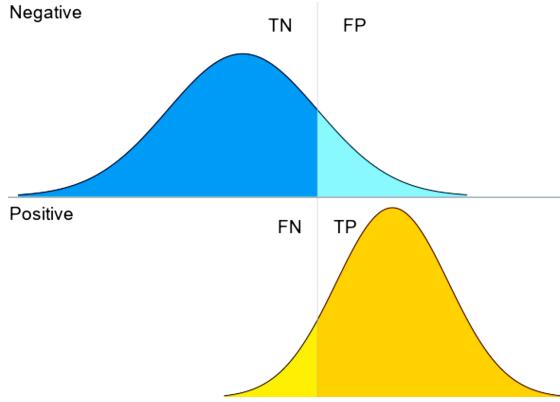


Figure 5.1: Relationship between threshold and classification. Figure obtained from [32]

Various metrics assess predictive model performance, such as:

- **Accuracy** is the ratio of correctly classified instances (TP) and (TN) to the total number of instances. It provides a general measure of the model's correctness.

$$\text{Accuracy } A = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Precision (P)** is the proportion of correctly predicted active cases (TP) to all instances predicted as active (TP + FP).

$$\text{Precision } P = \frac{TP}{TP + FP}$$

- **Recall (R) or Sensitivity or True Positive Rate (TPR)** is the proportion of correctly predicted active cases (TP) to all actual active cases (TP + FN).

$$\text{Recall } R = \frac{TP}{TP + FN}$$

- **F1 Score** is the harmonic mean of precision and recall, which balances the trade-off between false positives and false negatives.

$$F1 = \frac{2 \cdot P \cdot R}{P + R}$$

5.1. Binary Classification

- **True Negative Rate (TNR) or Specificity** is the proportion of correctly predicted inactive cases (TN) to all actual inactive cases (TN + FP).

$$TNR = \frac{TN}{TN + FP}$$

- **Receiver Operating Characteristic (ROC) Curve** is a graphical representation of the model's performance across different classification thresholds and plots the true positive rate (TPR) against the false positive rate (FPR), as exemplified for assay endpoint with *aeid* = 1120 in Figure 5.2. The Area Under the ROC Curve (AUC) is an indicator of the model's performance taking into account all possible classification thresholds, where a score of 1 represents a perfect model and 0.5 signifies a random no-skill model.

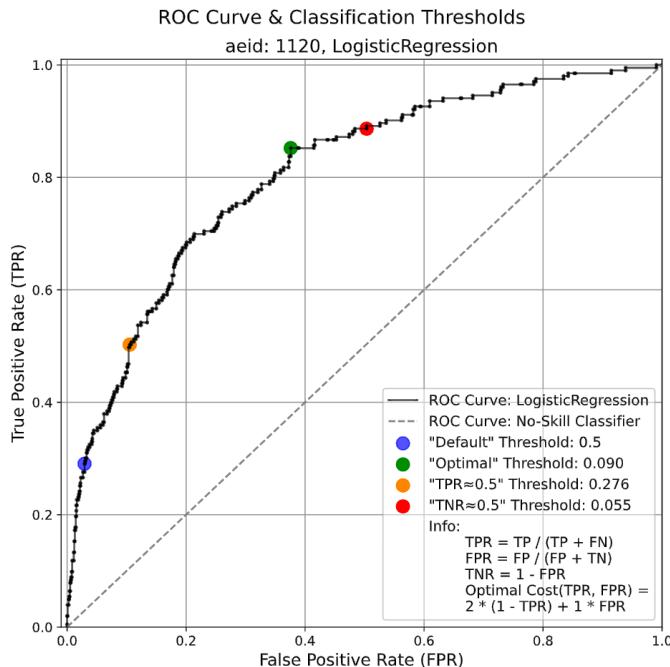


Figure 5.2: The Receiver Operating Characteristic (ROC) curve is presented for the LogisticRegression classifier in the case of assay endpoint with *aeid*: 1120. We make predictions for each model combination using four distinct classification thresholds, specifically: the default threshold of 0.5, the optimal threshold determined by the cost function weighting TPR twice as FPR (to value recall), TPR approximately equal to 0.5, and TNR approximately equal to 0.5. Figure 5.3 shows the confusion matrices for the four thresholds.

Imbalanced Data

A substantial proportion of the studied assay endpoints have an unequal distribution of active (positive) and inactive (negative) compounds. Typically, the negative class significantly outweighs the positive class, as depicted by an example in Figure 5.3.

5.1. Binary Classification

Confusion Matrices for 4 classification thresholds

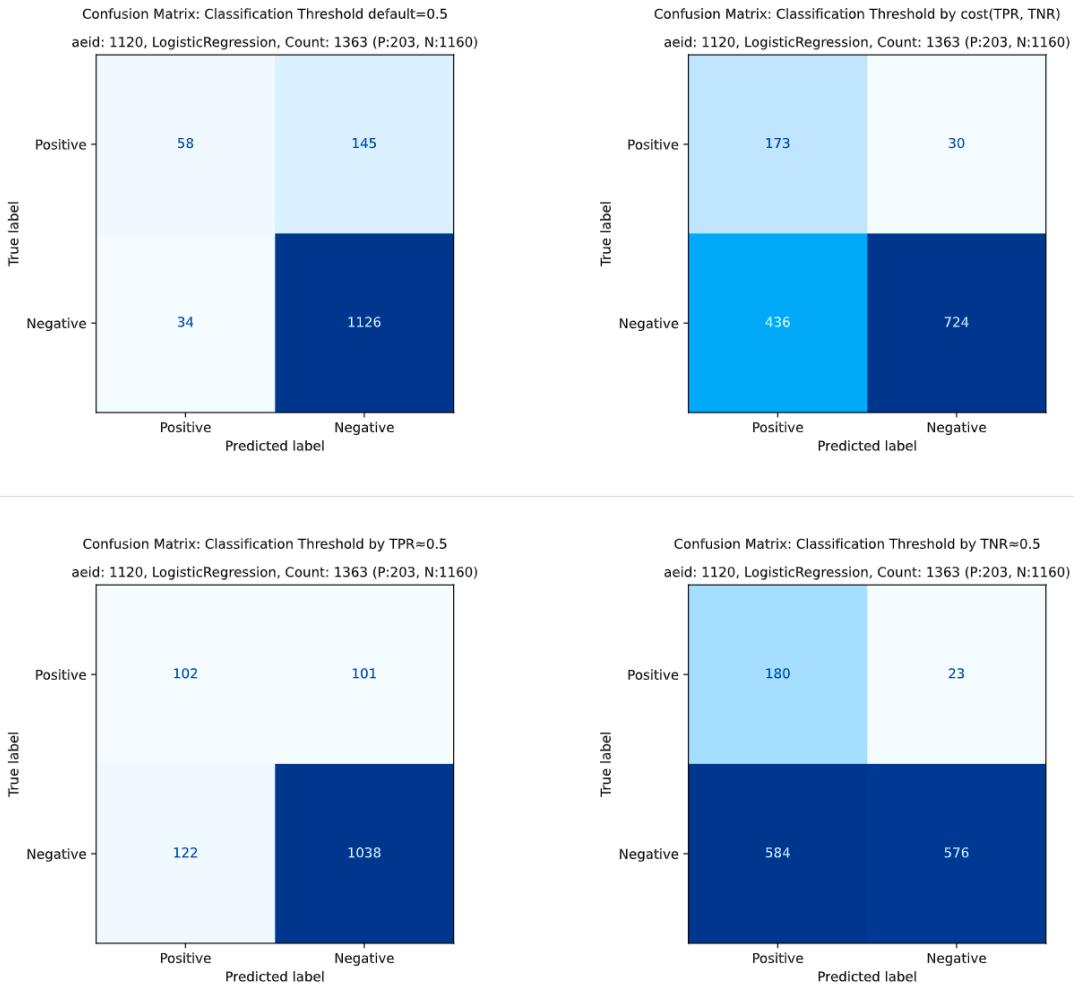


Figure 5.3: For assay endpoint with aeid: 1120, confusion matrices are shown for four different classification thresholds.

Such imbalanced datasets can result in skewed performance metrics, where the model may exhibit strong performance on the majority class while performing poorly on the minority class. To address imbalanced datasets, additional metrics such as macro averaged and weighted averaged metrics can be taken into account. In macro-averaging, individual class metrics are calculated and then equally weighted to compute an overall metric, disregarding class distribution. For instance, macro recall is the average of individual class recalls:

$$\text{Macro Recall} = \frac{R_{positive} + R_{negative}}{2}$$

In weighted averaging, individual class metrics are calculated and then weighted by the number of true instances in each class. For instance, weighted recall is the average of individual class recalls weighted by the number of true instances in each class. Similarly macro-averaged and weighted-averaged precision and F1 score can be calculated. The following two metric scores are not directly dependent on the threshold, and thus are well-suited for comparing models with imbalanced datasets:

- **Balanced Accuracy (BAC)** accounts for class imbalance by averaging the true positive rate (sensitivity) and true negative rate (specificity).

$$\text{Balanced Accuracy (BAC)} = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right)$$

- **PR-AUC (Precision-Recall Area Under the Curve)** quantifies the performance of the positive class by assessing the area under the precision-recall curve.

5.1.2 Performance

Performance results were generated for every combination of:

- Target variables (y):
 1. hitcall without cytotoxicity correction
 2. hitcall with cytotoxicity correction
- 345 Assay endpoints:
- Feature selection (FS) models:
 1. XGBoost
 2. RandomForest
- Estimator models:
 1. LogisticRegression
 2. MLPClassifier
 3. RandomForestClassifier
 4. SVM (Support Vector Machine)
 5. XGBoostClassifier
- Validation sets:
 1. the internal validation dataset
 2. MassBank validation set with fingerprints from chemical structure
 3. MassBank validation set with SIRIUS-predicted fingerprints

5.1. Binary Classification

- Classification thresholds:
 1. *default*: classification threshold of 0.5
 2. $TPR \simeq 0.5$: TPR approximately equal to 0.5. In other words, the threshold is chosen such that the models detect approximately 50% of the toxic compounds. This choice provides an estimate on the price to pay in terms of false positives, facilitating performance comparisons with other models.
 3. $TNR \simeq 0.5$: TNR approximately equal to 0.5.
 4. $\text{cost}(TPR, TNR)$: cost function weighting TPR twice as FPR: the threshold is chosen such that $\text{cost}(TPR, TNR) = 2 * (1 - TPR) + FPR$ is minimized, valuing recall twice as precision.
- Metrics on: 1. macro average, 2. weighted average, 3. positive and 4. negative class

Below, we present a selection of these outcomes in the form of figures and summary tables, showcasing the performance of binary classification models with the following fixed configurations:

- binarized hitcall without cytotoxicity correction as the target variable
- XGBoost as the feature selection model
- macro averaged metrics

With these settings, we evaluated the performance across all assay endpoints and estimators on the internal validation set the dual MassBank validation set. Note that reported balanced accuracy, the ROC-AUC and PR-AUC metric scores remain consistent across subsequent summaries, unaffected by the choice of classification threshold. The figures are presented in the following order:

Table 5.2: Figure References and Descriptions

Validation set	Classification Threshold	Figure
Internal	<i>default</i> = 0.5	Figure 5.4a
Internal	$\text{cost}(TPR, TNR) = 2 \cdot (1 - TPR) + FPR$	Figure 5.4b
Internal	$TPR \simeq 0.5$	Figure 5.5a
Internal	$TNR \simeq 0.5$	Figure 5.5b
MassBank from structure	<i>default</i> = 0.5	Figure 5.6a
MassBank SIRIUS-predicted	<i>default</i> = 0.5	Figure 5.6b
MassBank from structure	$\text{cost}(TPR, TNR) = 2 \cdot (1 - TPR) + FPR$	Figure 5.7a
MassBank SIRIUS-predicted	$\text{cost}(TPR, TNR) = 2 \cdot (1 - TPR) + FPR$	Figure 5.7b
MassBank from structure	$TPR \simeq 0.5$	Figure 5.8a
MassBank SIRIUS-predicted	$TPR \simeq 0.5$	Figure 5.8b
MassBank from structure	$TNR \simeq 0.5$	Figure 5.9a
MassBank SIRIUS-predicted	$TNR \simeq 0.5$	Figure 5.9b

5.1. Binary Classification

Internal Validation Set

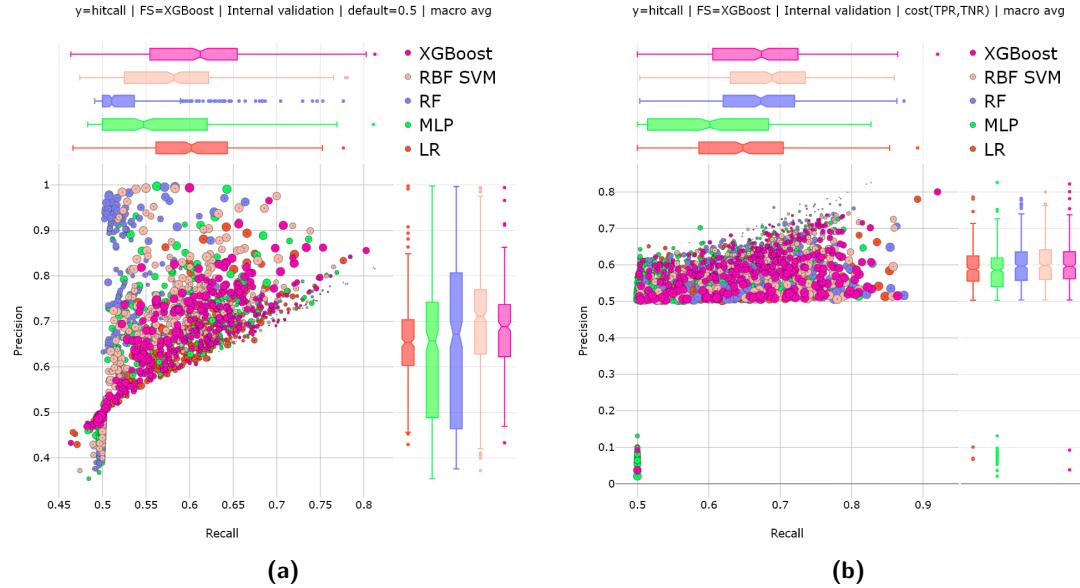


Figure 5.4: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints for the *internal validation set* (a) $\text{default} = 0.5$ (b) $\text{cost}(TPR, TNR) = 2 * (1 - TPR) + FPR$ threshold. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers). The tables below provide the median metrics for the estimators across all target assay endpoint models.

Table 5.3: Median Performance Metrics belonging to 5.4a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.653	0.602	0.617	0.827	0.602	0.711	0.367
MLP	0.657	0.547	0.547	0.844	0.547	0.678	0.339
RBF SVM	0.711	0.582	0.592	0.845	0.582	0.754	0.421
RF	0.671	0.511	0.491	0.846	0.511	0.744	0.392
XGBoost	0.688	0.612	0.632	0.837	0.612	0.741	0.417

Table 5.4: Median Performance Metrics belonging to 5.4b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.588	0.648	0.440	0.496	0.602	0.711	0.367
MLP	0.586	0.602	0.385	0.398	0.547	0.678	0.339
RBF SVM	0.598	0.689	0.510	0.559	0.582	0.754	0.421
RF	0.597	0.673	0.469	0.528	0.511	0.744	0.392
XGBoost	0.596	0.674	0.496	0.553	0.612	0.741	0.417

5.1. Binary Classification

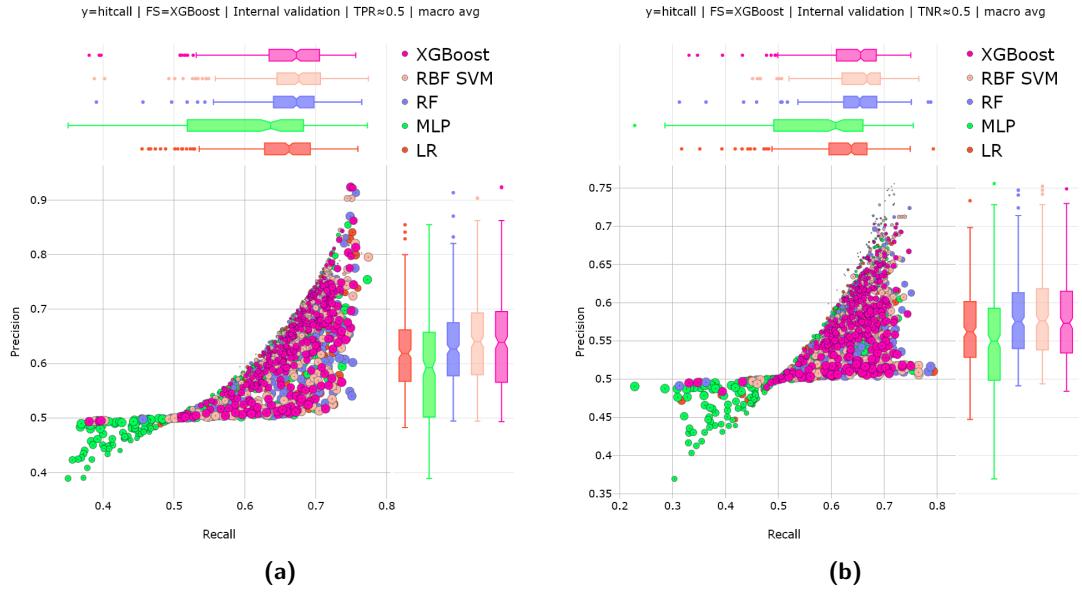


Figure 5.5: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints for the *internal validation set* (a) $TPR \approx 0.5$ (b) $TNR \approx 0.5$ classification threshold. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers). The tables below provide the median metrics for the estimators across all target assay endpoint models.

Table 5.5: Median Performance Metrics belonging to 5.5a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.619	0.662	0.632	0.746	0.602	0.711	0.367
MLP	0.592	0.636	0.597	0.702	0.547	0.678	0.339
RBF SVM	0.640	0.676	0.650	0.771	0.582	0.754	0.421
RF	0.627	0.673	0.639	0.765	0.511	0.744	0.392
XGBoost	0.639	0.673	0.649	0.763	0.612	0.741	0.417

Table 5.6: Median Performance Metrics belonging to 5.5b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.562	0.638	0.486	0.538	0.602	0.711	0.367
MLP	0.550	0.608	0.472	0.532	0.547	0.678	0.339
RBF SVM	0.576	0.667	0.498	0.547	0.582	0.754	0.421
RF	0.575	0.654	0.496	0.544	0.511	0.744	0.392
XGBoost	0.573	0.655	0.498	0.545	0.612	0.741	0.417

5.1. Binary Classification

MassBank Validation Set

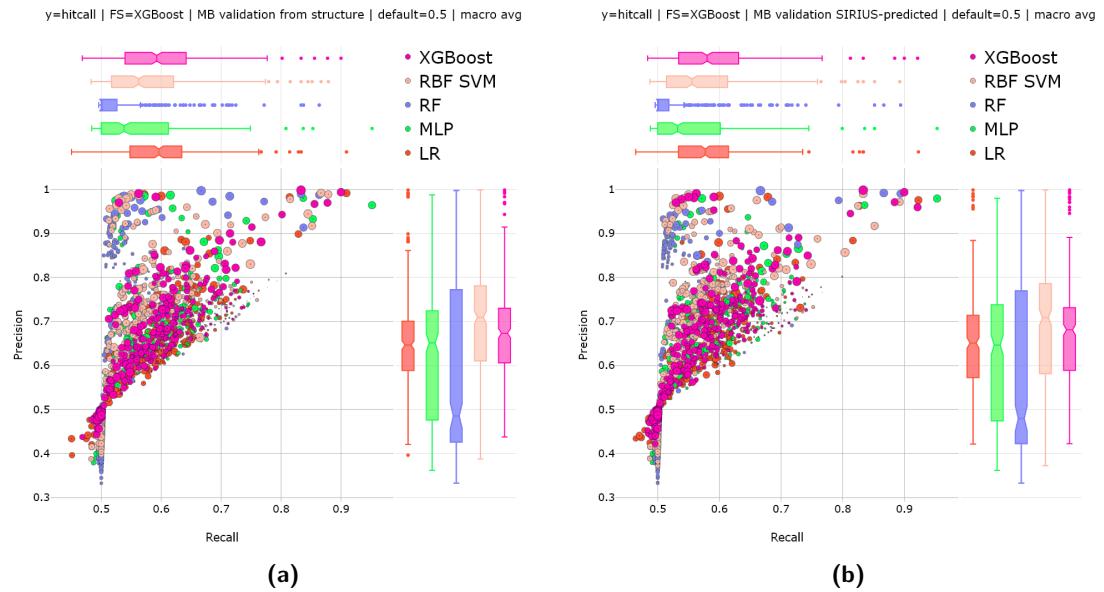


Figure 5.6: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints. MassBank validation with fingerprints from chemical structure (a) and SIRIUS-predicted fingerprints (b), $\text{default} = 0.5$ threshold, macro averaged metrics. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers).

Table 5.7: Median Performance Metrics belonging to 5.6a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.647	0.596	0.605	0.816	0.596	0.707	0.411
MLP	0.651	0.538	0.530	0.829	0.538	0.671	0.363
RBF SVM	0.709	0.562	0.570	0.835	0.562	0.738	0.455
RF	0.486	0.500	0.474	0.829	0.500	0.731	0.436
XGBoost	0.673	0.593	0.607	0.823	0.593	0.736	0.451

Table 5.8: Median Performance Metrics belonging to 5.6b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.651	0.577	0.587	0.820	0.577	0.693	0.385
MLP	0.646	0.532	0.519	0.829	0.532	0.668	0.355
RBF SVM	0.709	0.556	0.559	0.834	0.556	0.726	0.448
RF	0.480	0.500	0.471	0.829	0.500	0.727	0.433
XGBoost	0.681	0.580	0.594	0.829	0.580	0.721	0.441

5.1. Binary Classification

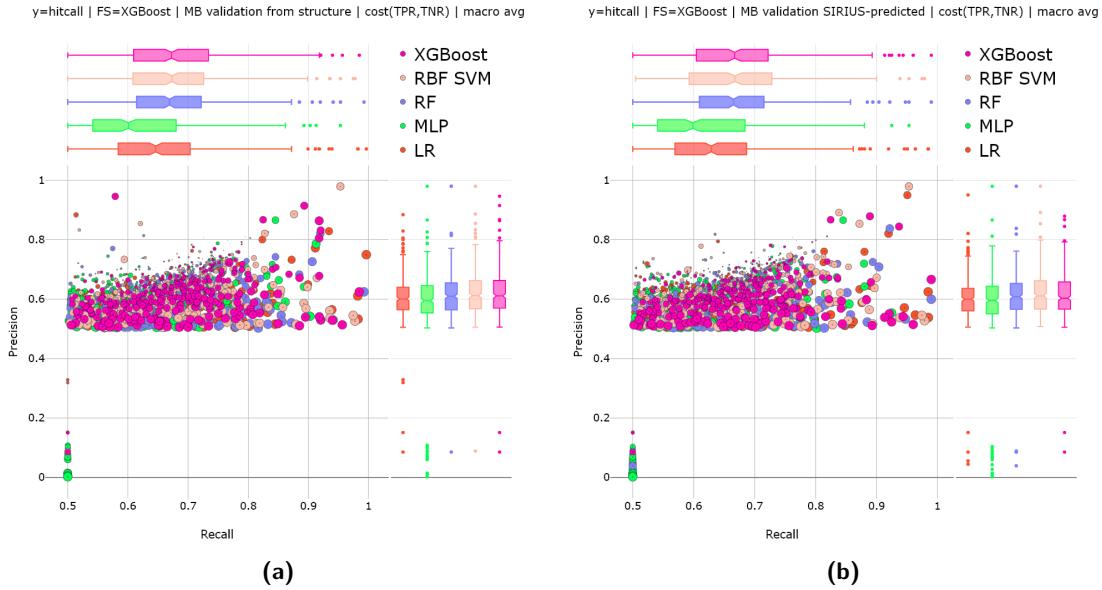


Figure 5.7: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints. MassBank validation with fingerprints from chemical structure (a) and SIRIUS-predicted fingerprints (b), $\text{cost}(TPR, TNR) = 2 * (1 - TPR) + FPR$ threshold, macro averaged metrics. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers). The tables below provide the median metrics for the estimators across all target assay endpoint models.

Table 5.9: Median Performance Metrics belonging to 5.7a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.601	0.646	0.470	0.503	0.596	0.707	0.411
MLP	0.595	0.601	0.386	0.407	0.538	0.671	0.363
RBF SVM	0.612	0.673	0.516	0.559	0.562	0.738	0.455
RF	0.610	0.669	0.507	0.548	0.500	0.731	0.436
XGBoost	0.611	0.673	0.511	0.563	0.593	0.736	0.451

Table 5.10: Median Performance Metrics belonging to 5.7b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.600	0.629	0.432	0.458	0.577	0.693	0.385
MLP	0.596	0.599	0.373	0.415	0.532	0.668	0.355
RBF SVM	0.611	0.668	0.493	0.559	0.556	0.726	0.448
RF	0.608	0.665	0.499	0.534	0.500	0.727	0.433
XGBoost	0.603	0.667	0.502	0.539	0.580	0.721	0.441

5.1. Binary Classification

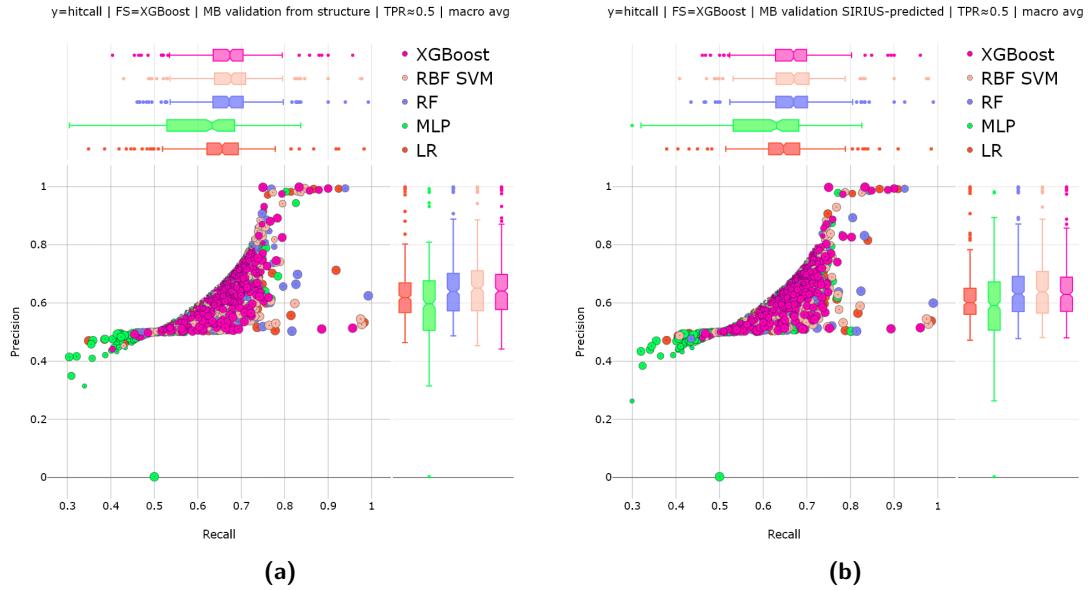


Figure 5.8: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints. MassBank validation with fingerprints from chemical structure (a) and SIRIUS-predicted fingerprints (b), $TPR \approx 0.5$ threshold, macro averaged metrics. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers). The tables below provide the median metrics for the estimators across all target assay endpoint models.

Table 5.11: Median Performance Metrics belonging to 5.8a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.618	0.656	0.626	0.727	0.596	0.707	0.411
MLP	0.597	0.633	0.605	0.683	0.538	0.671	0.363
RBF SVM	0.651	0.677	0.656	0.754	0.562	0.738	0.455
RF	0.639	0.672	0.647	0.750	0.500	0.731	0.436
XGBoost	0.641	0.674	0.649	0.754	0.593	0.736	0.451

Table 5.12: Median Performance Metrics belonging to 5.8b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.601	0.647	0.608	0.711	0.577	0.693	0.385
MLP	0.592	0.630	0.596	0.676	0.532	0.668	0.355
RBF SVM	0.638	0.671	0.644	0.749	0.556	0.726	0.448
RF	0.631	0.669	0.638	0.743	0.500	0.727	0.433
XGBoost	0.629	0.670	0.638	0.743	0.580	0.721	0.441

5.1. Binary Classification

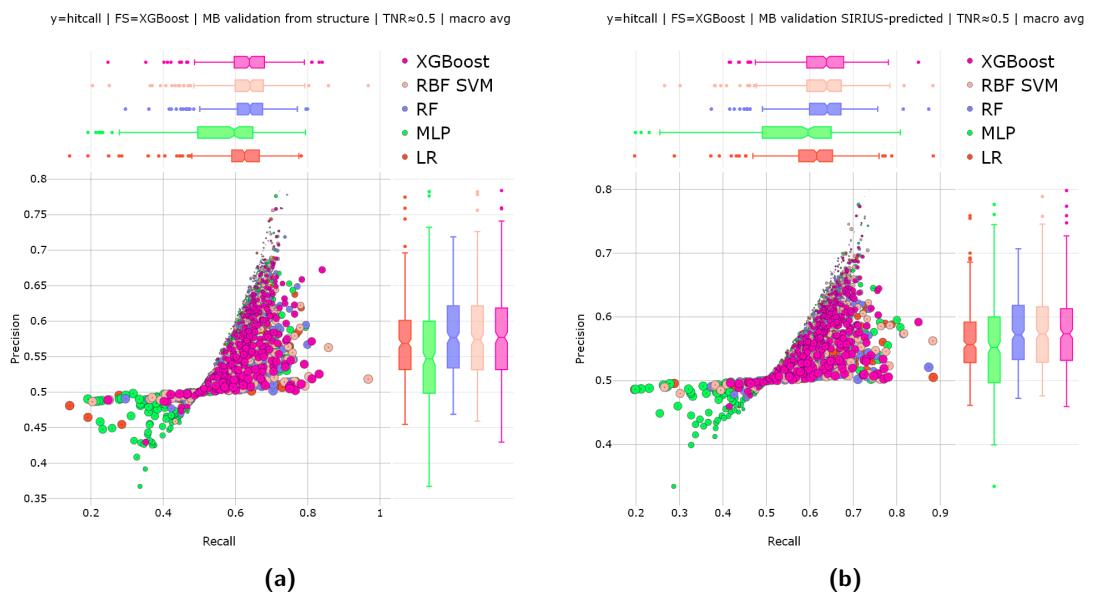


Figure 5.9: Comparison of precision and recall for five different estimators across a total of 345 evaluated assay endpoints. MassBank validation with fingerprints from chemical structure (a) and SIRIUS-predicted fingerprints (b), $TNR \approx 0.5$ threshold, macro averaged metrics. Larger marker size indicate a larger relative imbalance between the support for negative and positive compounds, normalized across all assay endpoints. The marginal boxplots illustrate the metric distribution across the target assay endpoint models (median, first quartile, third quartile, range-whiskers and outliers). The tables below provide the median metrics for the estimators across all target assay endpoint models.

Table 5.13: Median Performance Metrics belonging to 5.9a.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.569	0.625	0.509	0.552	0.596	0.707	0.411
MLP	0.547	0.596	0.478	0.537	0.538	0.671	0.363
RBF SVM	0.574	0.640	0.509	0.556	0.562	0.738	0.455
RF	0.577	0.641	0.514	0.555	0.500	0.731	0.436
XGBoost	0.577	0.638	0.509	0.553	0.593	0.736	0.451

Table 5.14: Median Performance Metrics belonging to 5.9b.

Estimator	Precision	Recall	F1	Acc.	Bal. Acc.	ROC-AUC	PR-AUC
LR	0.557	0.616	0.497	0.543	0.577	0.693	0.385
MLP	0.552	0.596	0.481	0.541	0.532	0.668	0.355
RBF SVM	0.573	0.639	0.506	0.555	0.556	0.726	0.448
RF	0.572	0.640	0.507	0.553	0.500	0.727	0.433
XGBoost	0.573	0.638	0.507	0.553	0.580	0.721	0.441

5.1.3 Feature Importance

An assessment of fingerprint feature importances enables to relate molecular substructures present in compounds to their toxicity. A high feature importance can signify that the presence or absence of the substructure is either linked to toxicity or non-toxicity. We gathered feature importance scores from the trained XGBoost and random forest classifier across all assay endpoints. For instance, we can visualize the top N features for each assay endpoints, which is shown for $N = 10$ in Figure 5.10.

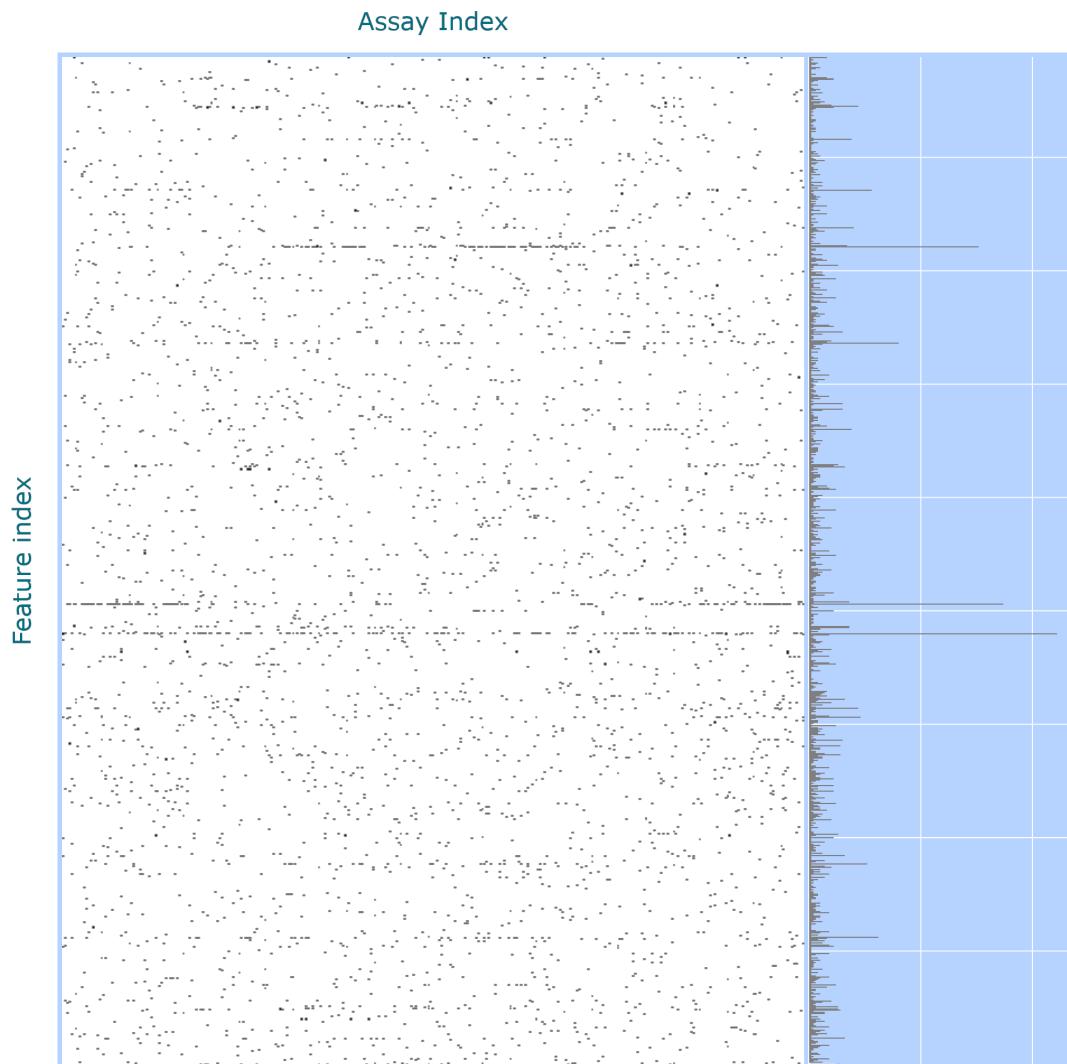


Figure 5.10: Top 10 Feature Importance Across All Target Assay Endpoints. Note that there are only a total of 889 features in the top 10 features for all assay endpoints combined. As a result, in order to make visualization easier, we reassigned the original feature indices to a new linear index.

Chapter 6

Discussion

When it comes to prioritizing compound identification in environmental samples through hazard assessment, the objective is to maximize the likelihood of detecting toxic substances while simultaneously minimizing the instances of false alarms which represent non-hazardous compounds misclassified as toxic.

Dealing with highly imbalanced datasets, coupled with a vast chemical space that is assumed to be sparsely represented in the training dataset, can present a challenging task for achieving robust predictive performance on unweighted class metrics. Nonetheless, the results demonstrate that the developed pipeline is capable of predicting toxicity based on molecular fingerprints derived from chemical structure. However, the model performances obtained are relatively modest, and we were unable to validate the reported balanced accuracy of 0.75 as found in [11].

In the case of the internal validation set, consider Figure 5.5a, where a macro-averaged metric with a particular threshold that fixes the true positive rate to approximately 0.5 is employed. The XGBoost classifier attains a median F1-score of 0.65 and balanced accuracy of 0.61 over 345 target assay endpoints in total. The performance of the multi-layer perceptron (MLP) was found to be disappointing and can be ascribed to the constraints imposed by the limited volume of data available for training a neural network.

The slight decline in model performance for the MassBank validation set based on chemical fingerprints from structure can be explained by the fact that the MassBank validation set's target variable distribution is not a perfect match with that of the training set, in contrast to the internal stratified validation set, which mirrors the training set's distribution. Also, a more in-depth analysis of the compounds within the MassBank validation set can shed light on any potential bias in the chemical space it covers. Additionally, it's important to note that the MassBank validation set is relatively small, which may contribute to increased variability in the model's performance. The model performances assessed on the MassBank validation set using SIRIUS-predicted fingerprints show a marginal dip compared to those relying on

chemical fingerprints derived from the molecular structure. This is an encouraging result, as it demonstrates that the developed pipeline is capable of predicting toxicity based on predicted fingerprints from spectral data.

Importantly, the results emphasize the crucial importance of the chosen classification threshold paired with the employed class metric (macro average, weighted average, positive, negative) in achieving this balance, especially when dealing with such imbalanced toxicity datasets. As an example, please refer to Figure 5.4a, where the sensitivity of the random forest classifier to the threshold is clearly evident. Notably, when employing the default threshold of 0.5, the random forest classifier exhibits a lower F1-score compared to the logistic regression classifier which tends to be more robust to threshold variations due to its inherent probabilistic nature. However, when considering alternative thresholds, the random forest classifier outperforms the logistic regression classifier in terms of F1-score. This sensitivity can be attributed to the highly imbalanced data.

In practical application of these models, selecting an ideal threshold depends on the laboratory's specific prioritization goals and resource constraints. This decision should take into account the sample size and the level of false alarms that can be tolerated, particularly if all positive predicted instances are pipelineed for an early-stage filtering process. In cases where prioritization aims to capture as many toxic compounds as possible, the threshold should be configured such that the true positive rate (TPR) takes precedence over the true negative rate (TNR). We explicitly examined a scenario (e.g., Figure 5.4a) in which TPR was given twice the weight of TNR, resulting in the detection of more toxic compounds but also incurring a higher number of false alarms. On the other hand, when laboratory testing resources are limited, setting a higher threshold may be preferable. While this mitigates the number of false alarms, it comes with the cost of a reduction in the number of detected toxic compounds.

Chapter 7

Conclusion

In summary, the performance results were explored across diverse combinations of models and metrics, providing a thorough understanding of the model's capabilities in various scenarios. This analysis allowed for a deeper comprehension of the strengths and weaknesses of the binary classification models employed in hazard assessment.

Bibliography

- [1] U. N. E. Programme, *Global Chemicals Outlook II - From Legacies to Innovative Solutions: Implementing the 2030 Agenda for Sustainable Development - Synthesis Report*, 2019. [Online]. Available: <https://wedocs.unep.org/20.500.11822/27651>.
- [2] C. A. Service, *Chemical Abstracts Service (CAS) is a division of the American Chemical Society*, Source of chemical information located in Columbus, Ohio, United States, <https://www.cas.org/support/documentation/cas-databases>, 2023.
- [3] R. Schwarzenbach *et al.*, “The Challenge of Micropollutants in Aquatic Systems,” *Science (New York, N.Y.)*, vol. 313, pp. 1072–7, Sep. 2006. doi: [10.1126/science.1127291](https://doi.org/10.1126/science.1127291).
- [4] E. Commission, D.-G. for Research, and Innovation, *European Green Deal - Research & innovation call*. Publications Office of the European Union, 2021. doi: [10.2777/33415](https://doi.org/10.2777/33415).
- [5] E. Commission, “EU Chemicals Strategy for Sustainability Towards a Toxic-Free Environment,” 2020, Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions - Chemicals Strategy for Sustainability Towards a Toxic-Free Environment. [Online]. Available: https://environment.ec.europa.eu/strategy/chemicals-strategy_en.
- [6] S. Tamara, M. A. den Boer, and A. J. R. Heck, “High-Resolution Native Mass Spectrometry,” *Chemical Reviews*, vol. 122, no. 8, pp. 7269–7326, 2022, PMID: 34415162. doi: [10.1021/acs.chemrev.1c00212](https://doi.org/10.1021/acs.chemrev.1c00212). eprint: <https://doi.org/10.1021/acs.chemrev.1c00212>. [Online]. Available: <https://doi.org/10.1021/acs.chemrev.1c00212>.
- [7] J. Hollender, E. L. Schymanski, H. P. Singer, and P. L. Ferguson, “Nontarget Screening with High Resolution Mass Spectrometry in the Environment: Ready to Go?” *Environmental Science & Technology*, vol. 51, no. 20, pp. 11 505–11 512, 2017, PMID: 28877430. doi: [10.1021/acs.est.7b02184](https://doi.org/10.1021/acs.est.7b02184). eprint: <https://doi.org/10.1021/acs.est.7b02184>.

Bibliography

- 1021/acs.est.7b02184. [Online]. Available: <https://doi.org/10.1021/acs.est.7b02184>.
- [8] P. Banerjee, A. O. Eckert, A. K. Schrey, and R. Preissner, "ProTox-II: a webserver for the prediction of toxicity of chemicals," *Nucleic Acids Research*, vol. 46, no. W1, W257–W263, Apr. 2018, ISSN: 0305-1048. doi: [10.1093/nar/gky318](https://doi.org/10.1093/nar/gky318). eprint: <https://academic.oup.com/nar/article-pdf/46/W1/W257/25110434/gky318.pdf>. [Online]. Available: <https://doi.org/10.1093/nar/gky318>.
- [9] N. H. G. R. I. Maggie Bartlett. "Chemical Genomics Robot." (2009), [Online]. Available: https://en.wikipedia.org/wiki/High-throughput_screening#/media/File:Chemical_Genomics_Robot.jpg.
- [10] J. Rudd. "High Throughput Screening - Accelerating Drug Discovery Efforts." (2017), [Online]. Available: <https://www.ddw-online.com/hts-a-strategy-for-drug-discovery-900-200008/>.
- [11] K. Arturi and J. Hollender, "Machine Learning-Based Hazard-Driven Prioritization of Features in Nontarget Screening of Environmental High-Resolution Mass Spectrometry Data," *Environmental Science & Technology*, vol. 0, no. 0, null, 0, PMID: 37279189. doi: [10.1021/acs.est.3c00304](https://doi.org/10.1021/acs.est.3c00304). eprint: <https://doi.org/10.1021/acs.est.3c00304>. [Online]. Available: <https://doi.org/10.1021/acs.est.3c00304>.
- [12] K. Dührkop *et al.*, "SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite structure information," *Nature methods*, vol. 16, no. 4, pp. 299–302, Apr. 2019, ISSN: 1548-7091. doi: [10.1038/s41592-019-0344-8](https://doi.org/10.1038/s41592-019-0344-8). [Online]. Available: https://research.aalto.fi/files/32997691/SCI_Duhrkop_Fleischauer_Sirius_4_Turning_tandem.pdf.
- [13] *MassBank: High Quality Mass Spectral Database*, <https://massbank.eu/MassBank/>, Accessed: 2023.
- [14] T. Janel, K. Takeuchi, and J. Bajorath, "Introducing a Chemically Intuitive Core-Substituent Fingerprint Designed to Explore Structural Requirements for Effective Similarity Searching and Machine Learning," *Molecules*, vol. 27, no. 7, 2022, ISSN: 1420-3049. doi: [10.3390/molecules27072331](https://doi.org/10.3390/molecules27072331). [Online]. Available: <https://www.mdpi.com/1420-3049/27/7/2331>.
- [15] S. M. Bell *et al.*, "In vitro to in vivo extrapolation for high throughput prioritization and decision making," *Toxicology in Vitro*, vol. 47, pp. 213–227, 2018, ISSN: 0887-2333. doi: <https://doi.org/10.1016/j.tiv.2017.11.016>. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0887233317303661>.
- [16] P. Nymark *et al.*, "Systematic Organization of COVID-19 Data Supported by the Adverse Outcome Pathway Framework," *Frontiers in Public Health*, vol. 9, May 2021. doi: [10.3389/fpubh.2021.638605](https://doi.org/10.3389/fpubh.2021.638605).

Bibliography

- [17] R. Judson *et al.*, “Editor’s Highlight: Analysis of the Effects of Cell Stress and Cytotoxicity on In Vitro Assay Activity Across a Diverse Chemical and Assay Space,” *Toxicological Sciences*, vol. 152, no. 2, pp. 323–339, May 2016, ISSN: 1096-6080. doi: [10.1093/toxsci/kfw092](https://doi.org/10.1093/toxsci/kfw092). eprint: <https://academic.oup.com/toxsci/article-pdf/152/2/323/26290632/kfw092.pdf>. [Online]. Available: <https://doi.org/10.1093/toxsci/kfw092>.
- [18] B. Escher, P. Neale, and F. Leusch, *Bioanalytical Tools in Water Quality Assessment*. IWA Publishing, Jun. 2021, ISBN: 9781789061987. doi: [10.2166/9781789061987](https://doi.org/10.2166/9781789061987). eprint: <https://iwaponline.com/book-pdf/899726/wio9781789061987.pdf>. [Online]. Available: <https://doi.org/10.2166/9781789061987>.
- [19] K. A. Fay *et al.*, “Differentiating Pathway-Specific From Nonspecific Effects in High-Throughput Toxicity Data: A Foundation for Prioritizing Adverse Outcome Pathway Development,” *Toxicological Sciences*, vol. 163, no. 2, pp. 500–515, Feb. 2018, ISSN: 1096-6080. doi: [10.1093/toxsci/kfy049](https://doi.org/10.1093/toxsci/kfy049). eprint: <https://academic.oup.com/toxsci/article-pdf/163/2/500/24935129/kfy049.pdf>. [Online]. Available: <https://doi.org/10.1093/toxsci/kfy049>.
- [20] C. N. Cavasotto and V. Scardino, “Machine Learning Toxicity Prediction: Latest Advances by Toxicity End Point,” *ACS Omega*, vol. 7, no. 51, pp. 47 536–47 546, 2022. doi: [10.1021/acsomega.2c05693](https://doi.org/10.1021/acsomega.2c05693). eprint: <https://doi.org/10.1021/acsomega.2c05693>. [Online]. Available: <https://doi.org/10.1021/acsomega.2c05693>.
- [21] A. M. Richard *et al.*, “The Tox21 10K Compound Library: Collaborative Chemistry Advancing Toxicology,” *Chemical Research in Toxicology*, vol. 34, no. 2, pp. 189–216, 2021, PMID: 33140634. doi: [10.1021/acs.chemrestox.0c00264](https://doi.org/10.1021/acs.chemrestox.0c00264). eprint: <https://doi.org/10.1021/acs.chemrestox.0c00264>. [Online]. Available: <https://doi.org/10.1021/acs.chemrestox.0c00264>.
- [22] F. Kretschmer, J. Seipp, M. Ludwig, G. W. Klau, and S. Böcker, “Small molecule machine learning: All models are wrong, some may not even be useful,” *bioRxiv*, 2023. doi: [10.1101/2023.03.27.534311](https://doi.org/10.1101/2023.03.27.534311). eprint: <https://www.biorxiv.org/content/early/2023/03/27/2023.03.27.534311.full.pdf>. [Online]. Available: <https://www.biorxiv.org/content/early/2023/03/27/2023.03.27.534311>.
- [23] P. Peets, W.-C. Wang, M. MacLeod, M. Breitholtz, J. W. Martin, and A. Kruve, “MS2Tox Machine Learning Tool for Predicting the Ecotoxicity of Unidentified Chemicals in Water by Nontarget LC-HRMS,” *Environmental Science & Technology*, vol. 56, no. 22, pp. 15 508–15 517, 2022, PMID: 36269851. doi: [10.1021/acs.est.2c02536](https://doi.org/10.1021/acs.est.2c02536). eprint: <https://doi.org/10.1021/acs.est.2c02536>. [Online]. Available: <https://doi.org/10.1021/acs.est.2c02536>.
- [24] A. J. Williams *et al.*, “The CompTox Chemistry Dashboard: a community data resource for environmental chemistry,” *Journal of Cheminformatics*, vol. 9, no. 1, p. 61, 2017. doi: [10.1186/s13321-017-0247-6](https://doi.org/10.1186/s13321-017-0247-6). [Online]. Available: <https://doi.org/10.1186/s13321-017-0247-6>.

Bibliography

- [25] L. Wu, R. Huang, I. V. Tetko, Z. Xia, J. Xu, and W. Tong, "Trade-off Predictivity and Explainability for Machine-Learning Powered Predictive Toxicology: An in-Depth Investigation with Tox21 Data Sets," *Chemical Research in Toxicology*, vol. 34, no. 2, pp. 541–549, 2021, PMID: 33513003. doi: [10.1021/acs.chemrestox.0c00373](https://doi.org/10.1021/acs.chemrestox.0c00373). eprint: <https://doi.org/10.1021/acs.chemrestox.0c00373>. [Online]. Available: <https://doi.org/10.1021/acs.chemrestox.0c00373>.
- [26] J. Phuong *et al.* "ToxCast Assay Annotation Data User Guide." (2014), [Online]. Available: https://www.epa.gov/sites/default/files/2015-08/documents/toxcast_annotation_data_users_guide_20141021.pdf.
- [27] T. Sheffield, J. Brown, S. Davidson, K. P. Friedman, and R. Judson, "tcplfit2: an R-language general purpose concentration-response modeling package," *Bioinformatics*, vol. 38, no. 4, pp. 1157–1158, Nov. 2021, ISSN: 1367-4803. doi: [10.1093/bioinformatics/btab779](https://doi.org/10.1093/bioinformatics/btab779). eprint: <https://academic.oup.com/bioinformatics/article-pdf/38/4/1157/50422999/btab779.pdf>. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btab779>.
- [28] C. for Computational Toxicology and U. E. Exposure, *tcpl v3.0 Data Processing*, R package vignette for the tcpl package v3.0, CRAN, 2023. [Online]. Available: https://cran.r-project.org/web/packages/tcpl/vignettes/Data_processing.html.
- [29] A. B. Daniel *et al.*, "Data curation to support toxicity assessments using the Integrated Chemical Environment," *Frontiers in Toxicology*, vol. 4, 2022, ISSN: 2673-3080. doi: [10.3389/ftox.2022.987848](https://doi.org/10.3389/ftox.2022.987848). [Online]. Available: <https://www.frontiersin.org/articles/10.3389/ftox.2022.987848>.
- [30] M. Feshuk *et al.*, "The ToxCast pipeline: updates to curve-fitting approaches and database structure," *Frontiers in Toxicology*, vol. 5, 2023, ISSN: 2673-3080. doi: [10.3389/ftox.2023.1275980](https://doi.org/10.3389/ftox.2023.1275980). [Online]. Available: <https://www.frontiersin.org/articles/10.3389/ftox.2023.1275980>.
- [31] E. D. Watt and R. S. Judson, "Uncertainty quantification in ToxCast high throughput screening," *PLOS ONE*, vol. 13, no. 7, pp. 1–23, Jul. 2018. doi: [10.1371/journal.pone.0196963](https://doi.org/10.1371/journal.pone.0196963). [Online]. Available: <https://doi.org/10.1371/journal.pone.0196963>.
- [32] R. Wicklin. "Visualization of a binary classification analysis." The DO Loop, February 24, 2020. (2020), [Online]. Available: <https://blogs.sas.com/content/iml/2020/02/24/binary-classification-viz.html>.

Appendix A

Appendix

A.1 Variability in the Tested Concentration Across Assay Endpoints and Compounds

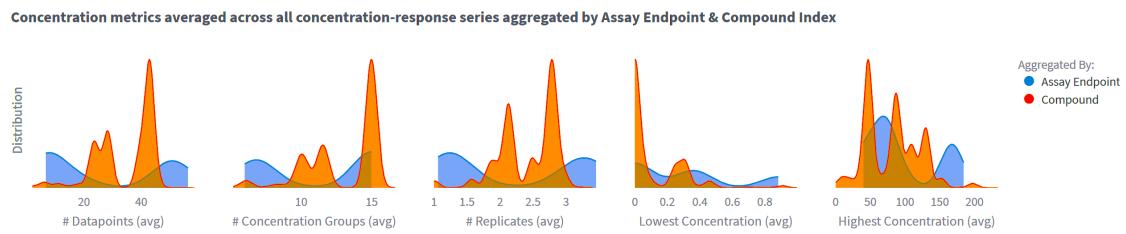


Figure A.1: Concentration metrics averaged across all concentration-response series aggregated by assay endpoint (blue) and compound (orange). E.g., the first chart shows the distribution on the average number of datapoints across all assay endpoint $a_i \in A$ with $\frac{1}{|A|} \sum_j n_{\text{datapoints}_{i,j}}$ and across all compounds $c_j \in C$ with $\frac{1}{|C|} \sum_i n_{\text{datapoints}_{i,j}}$. Similarly, the process is repeated for the other metrics: $n_{\text{groups}_{i,j}}$, $n_{\text{replicates}_{i,j}}$, $\min_{\text{conc}_{i,j}}$, and $\max_{\text{conc}_{i,j}}$.

A.2 ToxCast Assay Sources

Table A.1: Assay source names and long names

Assay source name	Assay source long name
ACEA	ACEA Biosciences
APR	Apredica
ATG	Attagene
BSK	Bioseek
NVS	Novascreen
OT	Odyssey Thera
TOX21	Tox21/NCGC
CEETOX	Ceetox/OpAns
LTEA	LifeTech/Expression Analysis
VALA	VALA Sciences
CLD	CellzDirect
CCTE_PADILLA	CCTE Padilla Lab
TANGUAY	Tanguay Lab
STM	Stemina Biomarker Discovery
ARUNA	ArunA Biomedical
CCTE	CCTE Labs
CCTE_SHAFER	CCTE Shafer Lab
CPHEA_STOKER	CPHEA Stoker and Laws Labs
CCTE_GLTED	CCTE Great Lakes Toxicology and Ecology Division
UPITT	University of Pittsburgh Johnston Lab
UKN	University of Konstanz
ERF	Eurofins
TAMU	Texas A&M University
IUF	Leibniz Research Institute for Environmental Medicine
CCTE_MUNDY	CCTE Mundy Lab
UTOR	University of Toronto, Peng Laboratory

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

ENHANCING TOXICITY PREDICTION OF MLINVITROTOX: PRIORITIZING UNIDENTIFIED COMPOUNDS IN ENVIRONMENTAL SAMPLES BASED ON HAZARD ASSESSMENT

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

BOSSHARD

First name(s):

ROBIN

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

16.10.2023

Signature(s)

R. Bosshard

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.